



2018

Self-Consistency In Sequential Decision-Making

Long Luu

University of Pennsylvania, thelong20.4@gmail.com

Follow this and additional works at: <https://repository.upenn.edu/edissertations>



Part of the [Neuroscience and Neurobiology Commons](#), and the [Psychology Commons](#)

Recommended Citation

Luu, Long, "Self-Consistency In Sequential Decision-Making" (2018). *Publicly Accessible Penn Dissertations*. 3152.
<https://repository.upenn.edu/edissertations/3152>

This paper is posted at ScholarlyCommons. <https://repository.upenn.edu/edissertations/3152>
For more information, please contact repository@pobox.upenn.edu.

Self-Consistency In Sequential Decision-Making

Abstract

Human decisions are rarely made in isolation. We typically have to make a sequence of decisions to reach a goal. Studies in economics and cognitive psychology have shown that making a decision may result in several biases in subsequent judgments. Similar biases have also recently been found in human percepts of low-level stimuli such as motion direction. What lacking is a principled framework that can account for several sequential dependencies between judgments. Towards that goal, in my thesis, I propose and experimentally test a self-consistent Bayesian observer model that assumes humans maintain self-consistency along the inference process. In Chapter 2, I first demonstrate that after having made a categorical decision on stimulus orientation, subjects' estimate of the stimulus is systematically biased away from the decision boundary. Two additional experiments suggest that the bias occurs because subjects treat their first decision as a fact and use that to constrain the subsequent estimation. Model fit to the data in my experiments and data in previous studies show that the self-consistent Bayesian model can quantitatively account for human behaviors in a wide range of experimental settings. In Chapter 3, using the same decision-estimation tasks, I probed the post-decision sensory representation by providing feedback on the categorical decision. I found that subjects' sensory representation is kept intact and the self-consistency is implemented by conditioning the prior distribution on the categorical decision. The results also suggest another interesting form of self-consistency when subjects' decision was incorrect: they reconstructed the sensory measurement to make it consistent with the given feedback. In Chapter 4, I found that the choice-induced bias also occurs in human judgment of number. The bias is similar for both non-symbolic (cloud of dots) and symbolic (sequence of Arabic numerals) forms of number. Finally, I propose in the general discussion how the self-consistent Bayesian framework may account for other biases in sequential decision-making such as the halo effect and sunk-cost fallacy.

Degree Type

Dissertation

Degree Name

Doctor of Philosophy (PhD)

Graduate Group

Psychology

First Advisor

Alan A. Stocker

Keywords

Bayesian decision theory, Decision-making, Perception, Psychophysics, Self-consistency, Sequential judgment

Subject Categories

Neuroscience and Neurobiology | Psychology | Social and Behavioral Sciences

SELF-CONSISTENCY IN SEQUENTIAL DECISION-MAKING

Long Luu

A DISSERTATION

in

Psychology

Presented to the Faculties of the University of Pennsylvania

in

Partial Fulfillment of the Requirements for the

Degree of Doctor of Philosophy

2018

Supervisor of Dissertation

Alan A. Stocker, Associate Professor of Psychology

Graduate Group Chairperson

Sara Jaffee, Professor of Psychology

Dissertation Committee

David H. Brainard, Professor of Psychology

Joshua I. Gold, Professor of Neuroscience

Joseph W. Kable, Associate Professor of Psychology

Dedicated to my parents and my wife
In loving memory of my third uncle Vuong

Acknowledgment

I would like to express my gratitude to many people who have helped to bring this thesis to fruition.

First I want to thank my advisor, Alan Stocker, for his dedicated mentoring. I have learned from him how to become a proper scientist which ranges from doing solid science and efficiently communicating my works to the scientific community to developing a scientific career. In every aspect of the scientific enterprise, Alan has shown me the high standards I should always aspire to.

I am also deeply grateful to my thesis committee, David Brainard, Josh Gold and Joe Kable for their relentless support and helpful guidance during my PhD years. I want to especially thank David for introducing me to the beauty and fascination of vision science and psychophysics. I want to thank Josh for always giving constructive feedback and helping me see the problem from a different perspective. I am grateful to Joe for his willingness to help when I need it.

I am grateful to many people in Psychology Department. First I would like to thank the administrative staffs, especially Yuni Thornton, Laurel Sweeney and Jessica Marcus who have helped me a lot with the administrative works. I also want to thank many faculty members in psychology who have introduced me to the exciting field of psychology through proseminar courses.

During the years at Penn, I have received generous mentoring and enjoyed friendship with many fellow students and postdocs, especially those in perception group. Manuel Spitschan provided me with detailed and valuable guidance in performing psychophysical experiment. Xuexin Wei always gave me interesting insights in our frequent conversation both inside and outside the lab. I want to thank Ana Radonjic, Nicolas Cottaris, Toni Sareela, Matjaz Jogan, Pedro Ortega, Alan He, Mohammad Rostami, Daniel Pak, Benjamin Chin, Cheng Qiu, Jiang Mao, Lingqi Zhang, Noam Roth, David White, Michael Barnett, Yunshu Fan and Takahiro Doi for many interesting discussions and conversations over science and life.

Last but not least, I want to thank my parents, my uncle Vuong and my wife for their relentless and unconditional support for my professional career no matter what path I have chosen or how esoteric my choice appears to them.

ABSTRACT

SELF-CONSISTENCY IN SEQUENTIAL DECISION-MAKING

Long Luu

Alan A. Stocker

Human decisions are rarely made in isolation. We typically have to make a sequence of decisions to reach a goal. Studies in economics and cognitive psychology have shown that making a decision may result in several biases in subsequent judgments. Similar biases have also recently been found in human percepts of low-level stimuli such as motion direction. What lacking is a principled framework that can account for several sequential dependencies between judgments. Towards that goal, in my thesis, I propose and experimentally test a self-consistent Bayesian observer model that assumes humans maintain self-consistency along the inference process. In Chapter 2, I first demonstrate that after having made a categorical decision on stimulus orientation, subjects' estimate of the stimulus is systematically biased away from the decision boundary. Two additional experiments suggest that the bias occurs because subjects treat their first decision as a fact and use that to constrain the subsequent estimation. Model fit to the data in my experiments and data in previous studies show that the self-consistent Bayesian model can quantitatively account for human behaviors in a wide range of experimental settings. In Chapter 3, using the same decision-estimation tasks, I probed the post-decision sensory representation by providing feedback on the categorical decision. I found that subjects' sensory representation is kept intact and the self-consistency is implemented by conditioning the prior distribution on the categorical decision. The results also suggest another interesting form of self-consistency when subjects' decision was incorrect: they reconstructed the sensory measurement to make it consistent with the given feedback. In Chapter 4, I found that the choice-induced bias also occurs in human judgment of number. The bias is similar for both non-symbolic (cloud of dots) and symbolic (sequence of Arabic numerals) forms of number. Finally, I propose in the general discussion how the self-consistent Bayesian framework may account for other biases in sequential decision-making such as the halo effect and sunk-cost fallacy.

TABLE OF CONTENTS

Abstract	v
List of figures	x
1 Introduction	1
1.1 Sequential dependency in cognitive judgments	3
1.2 Sequential dependency in perceptual judgments	5
1.2.1 Single-judgment task	5
1.2.2 Multiple-judgment task	7
1.3 Bayesian decision theory	8
1.4 Supplemental-psychophysics methods	12
2 Self-consistent inference	15
2.1 Introduction	15
2.2 Replicating choice-induced bias	16
2.3 Self-consistent Bayesian observer model	18
2.4 Validating self-consistent Bayesian observer model	23
2.4.1 Key features of the model	23
2.4.2 Inconsistent trials	28
2.4.3 Explaining existing experimental data	30
2.5 Self-consistency despite memory degradation	32
2.6 Alternative interpretations	33
2.7 Discussion	36
2.8 Method and supplementary material	37
3 Post-decision sensory representation	57
3.1 Introduction	57

3.2	Probing post-decision sensory representation	58
3.2.1	Experimental test and alternative models	58
3.2.2	Some sensory information is still preserved	61
3.3	Further test of alternative models	64
3.3.1	Imposing asymmetric prior	64
3.3.2	Subjects maintain full sensory representation	65
3.4	Discussion	68
3.5	Methods	70
4	Self-consistency in number judgment	79
4.1	Introduction	79
4.2	Dot array stimulus	80
4.2.1	Choice-induced bias in perceived number of dots	80
4.2.2	Explaining the bias with self-consistent model	81
4.3	Symbolic number stimulus	84
4.3.1	Probabilistic inference over a sequence of numbers	84
4.3.2	Trial-by-trial prediction of observer models	84
4.3.3	Replicating the result with different task instruction	87
4.4	Discussion	89
4.5	Method and supplementary material	92
5	General discussion	100
5.1	Summary of thesis contribution	100
5.2	Testable predictions of self-consistency principle	101
5.2.1	Sequential judgments on multiple attributes	101
5.2.2	Incorporating additional evidence after a preliminary decision	107
5.3	Self-consistent inference	109
5.3.1	Benefits of self-consistency	109
5.3.2	Conditioning likelihood and loss function	112

5.4	General implications for decision-making study	114
	Bibliography	117

List of figures

Figure 1.1: Example to illustrate Bayesian decision theory	14
Figure 2.1: Post-decision biases in a perceptual task sequence	17
Figure 2.2: Bayesian observer models for the perceptual task sequence	19
Figure 2.3: The self-consistent Bayesian observer model	21
Figure 2.4: Experiment 1: Data and model fits for individual subjects	22
Figure 2.5: Effect of the stimulus prior	24
Figure 2.6: Self-made vs. given category assignment	25
Figure 2.7: Experiments 2 and 3: Joint fit to data for individual subjects	27
Figure 2.8: Inconsistent trials are due to lapses and motor noise	29
Figure 2.9: Model fits for experimental data by Zamboni et al. (2016)	31
Figure 2.10: Maintaining self-consistency in the face of working memory noise	34
Figure 2.11: Full distributions of individual subjects' estimates in Experiment 1	47
Figure 2.12: Measured motor noise of individual subjects	48
Figure 2.13: Histogram plots of the orientation estimates together with the model fit for Experiment 1 (combined subject)	49
Figure 2.14: Goodness-of-fits for Experiment 1	50
Figure 2.15: Full distributions of individual subjects' estimates in Experiment 2	51
Figure 2.16: Full distributions of individual subjects' estimates in Experiment 3	52
Figure 2.17: Goodness-of-fits for Experiments 2 and 3	53
Figure 2.18: Zamboni et al. (2016) data (Experiment 1, combined subject) and fit with the self-consistent observer model	54
Figure 2.19: Zamboni et al. (2016) data (Experiment 2, combined subject) and fit with the self-consistent observer model	55
Figure 2.20: Relative log-likelihoods of model fits for Zamboni et al. (2016) data	56

Figure 3.1: Alternative hypotheses of post-decision sensory representation. . . .	59
Figure 3.2: Prediction of alternative hypotheses for incorrect trials.	60
Figure 3.3: Self-consistent model fit to correct trials.	62
Figure 3.4: Data and model prediction for incorrect trials.	63
Figure 3.5: Performance of alternative models in Experiment 1.	65
Figure 3.6: Dissociating the models by imposing an asymmetric prior.	66
Figure 3.7: Model prediction and data in Experiment 2.	67
Figure 3.8: Performance of alternative models in Experiment 2.	69
Figure 4.1: Experiment 1 procedure	80
Figure 4.2: Data and model fit in Experiment 1	82
Figure 4.3: Experiment 2 procedure and result	85
Figure 4.4: Comparing alternative models	86
Figure 4.5: Instruction for Experiment 3	88
Figure 4.6: Experiment 3 result	90
Figure 4.7: Histogram of subjects' estimates and model fits in Experiment 1 . .	98
Figure 4.8: Correlation of sample position and decision	99
Figure 5.1: Rational Bayesian observer of multi-attribute judgment	104
Figure 5.2: Self-consistent observer of multi-attribute judgment	105
Figure 5.3: Experimental test of halo effect at low-level perception	106
Figure 5.4: Modeling sunk cost fallacy	110
Figure 5.5: Advantage of self-consistent Bayesian observer	111
Figure 5.6: Conditioning the likelihood	113

Chapter 1

Introduction

As the field of machine intelligence advances at a blistering pace, we are now used to hearing breaking news about machine beating human experts on many areas once thought to be exclusively for the human. That ranges from AlphaGo beating Go champions (Silver et al., 2016, 2017), machine learning algorithms making medical decisions (Obermeyer and Emanuel, 2016; Kononenko, 2001) to autonomous robots fighting in wars (Singer, 2010; Ayoub and Payne, 2016). It leaves many of us to wonder what will be left for human intelligence. A common response is that we are still far from building a machine that makes a medical diagnosis in the morning, then drives to a friend's home to play a chess game in the afternoon and is ready to get in war at any time if necessary. In other words, the best machines excel humans at a specific area yet fall far behind an average person in all other areas (Lake et al., 2017; Kasparov, 2017).

Among areas that machines are still behind human in many regards is perception which involves judgments about basic properties of the world. Deep Blue computer that defeated Gary Kasparov cannot tell whether the sky is blue or how deep the ocean is. The ease with which we make such perceptual decisions conceals the complexity of the underlying processes. For instance, driving appears to be a fairly basic skill for most people, yet it requires a tremendous amount of complicated decisions many of which involve inference about perceptual features of the world (e.g. what are the surrounding objects and how far away or how fast are they?). The challenges are more obvious with the development of state-of-the-art technologies such as self-driving cars or humanoid robots. In fact, research on perception dates back to as early as the 19th century when the pioneers like Gustav Fechner and Hermann von Helmholtz conceived the lawful relations between the mental process and the physical world (Fechner, 1860; Helmholtz, 1867; Weber, 1867). An intensive body of

works has followed those early attempts and greatly advanced our understanding of how humans perceive the world (Cornsweet, 2012; Blakemore, 1993; Wandell, 1995; Green and Swets, 1966; Wolfe et al., 2006). In these studies, human subjects typically make a single judgment on a simple stimulus in each trial (see Supplemental-psychophysics methods for a brief overview of experimental methods). Crucially, the dominant theoretical frameworks in perception such as linear system theory (Wandell, 1995), signal detection theory (Macmillan and Creelman, 2004) and Bayesian decision theory (Kersten and Schrater, 2002; Maloney, 2002) mostly assume that perceptual judgments are made independent of each other. As a result, a perceptual decision should not affect subsequent judgments.

However, several lines of research in cognition and perception suggest that there exist many forms of sequential dependency between decisions. In cognitive psychology and economics, a rich body of works found evidence that a decision may cause systematic biases in subsequent judgments (Abelson, 1968; Kahneman, 2011). For instance, when humans made a sequence of judgments on the same stimulus features, the later judgment tend to be distorted to be consistent with the preceding decision (Brehm, 1956). Several other decisional biases induced by preceding judgments have also been well documented such as confirmation bias (Nickerson, 1998), halo effect (Nisbett and Wilson, 1977) and sunk-cost fallacy (Arkes and Blumer, 1985).

Several perception studies have found that perceptual decisions across trials exhibit specific patterns such as repetition or alternation despite the fact that the experimental stimuli were randomized (Fernberger, 1920; Gold et al., 2008; Fründ et al., 2014). An important difference between these studies and the cognitive literature mentioned above is that subjects only made one decision in each trial, which makes it difficult to have a direct comparison. More recently, a small body of works in perception had multiple judgments in each trial and show biases that are similar to the cognitive biases (Jazayeri and Movshon, 2007; Zamboni et al., 2016; Bronfman et al., 2015). As an example, a study by Jazayeri and Movshon (2007) shows that after having made a decision on motion direction of a dot cloud, subjects' percept of the stimulus motion direction was systematically biased away

from the decision boundary.

Although there exist many ad hoc theories that can explain each of these bias effects, what we lack is a unifying model framework that may account for many of the findings and provide fruitful predictions for future research. Towards that goal, my thesis extends a self-conditioned Bayesian observer model (Stocker and Simoncelli, 2007) and demonstrates through several experiments how the model can quantitatively account for sequential dependency in human perceptual decisions across various experimental settings.

In the rest of the chapter, I provide the necessary background for the readers that are not familiar with the relevant literature. Specifically, I will review key studies that show several forms of sequential dependency between decisions in cognition and perception. Then I will review the standard framework of Bayesian decision theory as applied in perception.

1.1. Sequential dependency in cognitive judgments

In 1956, the psychologist Leon Festinger published a book to describe the behavior of members in a religious cult that believes the world would be destroyed by a flood on December 21, 1954 and the believers would be saved by a flying saucer (Festinger et al., 1956). Interestingly, after they found out that the end of the world did not happen, they readily accepted the leader's argument that the world is saved by their extreme piety and even doubled down on their belief and fervently proselytized for it. To explain the phenomenon, Festinger developed a theory of **cognitive dissonance** which posits that people dislike holding contradictory beliefs in mind and tend to modify them to maintain a sense of consistency (Festinger, 1957). In the religious cult example, the members' belief that the world would be destroyed is highly dissonant with the belief that the event did not happen. Importantly, they cannot change either of them because they had given up everything (jobs, spouses, possessions) to follow the cult. Therefore, they easily fell prey to demagoguery claims and actively advanced those arguments that may help in reconciling these dissonant beliefs.

Following on Festinger's pioneering works, a large body of experimental studies has found evidence that is consistent with cognitive dissonance theory (Abelson, 1968; Festinger

and Carlsmith, 1959; Bem, 1972; Knox and Inkster, 1968; Younger et al., 1977; Frenkel and Doob, 1976; Cohen and Goldberg, 1970). In a study by Brehm (1956) subjects were asked to rate the desirability of eight items (e.g. coffee-maker, toaster, stopwatch). Then they had to choose between two equally desirable options (e.g. coffee-maker vs. toaster) and rate the items again. In accordance with cognitive dissonance theory, subjects' ratings of the items after a difficult decision change according to their choice, that is, the value of chosen item increases whereas the value of rejected item decreases. This line of works has been revived recently (Lee and Schwarz, 2010) and researchers recast it as choice-induced preference (Egan et al., 2010; Sharot et al., 2010, 2012) or preservation of coherence/consistency (Simon et al., 2004; Rieffer et al., 2017). By using choice-blindness paradigm (Johansson et al., 2005), the studies on choice-induced preference confirm that cognitive dissonance effect is not due to statistical artifacts as pointed out by Chen and Risen (2010).

As reviewed above, cognitive dissonance studies demonstrate sequential dependency between judgments that subjects make on the same feature using the same evidence. However, the practical situation that a decision maker encounters often involves judgments on different attributes (e.g. intelligence, experience, and integrity of a job candidate). If the judgments are made sequentially, the order in which the attributes are considered may significantly affect the judgment results as demonstrated in studies on **halo effect** (Nisbett and Wilson, 1977; Thorndike, 1920; Sigall and Ostrove, 1975; Leuthesser et al., 1995). Another typical real-life scenario is when we cannot acquire all relevant information at once. Hence, we first make a preliminary decision and subsequently incorporate any additional information presented to us. For example, a company decides to invest in a promising project with an initial estimated budget of \$100 million dollars. Halfway through the project, all money is spent and it is estimated that another \$100 million dollars are needed to finish it. Moreover, the project is not as promising as it was and the same amount of money can be spent on another more profitable project. What should the company do? Anecdotal observation and empirical evidence have shown that the company is more likely to continue the existing project. Noticeably, if it had not invested in the existing project, it would more

likely invest in the new project. This phenomenon is named **sunk-cost fallacy** (Arkes and Blumer, 1985; Staw, 1976). A related "irrational" behavior, confirmation bias, was also well documented in several studies which suggest that subjects tend to search for information supporting their previous decision (Nickerson, 1998; Mynatt et al., 1977; Jonas et al., 2001). Those scenarios will be discussed more extensively in section 5.2.

1.2. Sequential dependency in perceptual judgments

1.2.1. Single-judgment task

Even in the standard psychophysics experiments that include a single decision in a trial, researchers have long noticed some dependency in subjects' decisions across trials. As early as 1920, Fernberger (1920) found that when subjects made a simple discrimination on a pair of weights ("lighter", "heavier" or "equal"), their decision in the current trial was repulsed away from the decision in the previous trial (e.g. they tended to respond "lighter" when the response had been "heavier" in the preceding trial). Several subsequent studies using a wide range of stimuli and task designs found mixed results. In some studies, subjects tended to avoid the previous response (Arons and Irwin, 1932; Preston, 1936; Collier, 1951) whereas other studies show that subjects tended to repeat the previous response (Senders and Sowards, 1952; Senders, 1953; Verplanck et al., 1952; Day, 1956; McGill, 1957). Several studies have tried to determine the factors that influence this dependency such as training (Verplanck et al., 1953), intertrial interval, experimental duration (Collier, 1954; Day, 1957) and previous stimulus identity (Verplanck and Blough, 1958; Collier and Verplanck, 1958). More recent studies replicated the sequential dependency and demonstrated that the effect cannot be explained by motor bias, sensory adaptation or attentional effect (Akaishi et al., 2014; Fründ et al., 2014; Braun et al., 2018). When feedback was given in every trial, subjects' decision is dependent on the feedback in the previous trial: repetition if correct and alternation if incorrect (win-stay, lose-switch) (Busse et al., 2011; Abrahamyan et al., 2016; Urai et al., 2017).

How can we explain this type of sequential dependency? Some proposed that the dependency comes from subjects' response bias, that is, subjects' responses tend to form

a pattern even when the stimulus is totally ambiguous (Tune, 1964; Goodfellow, 1938). Other researchers approached the problem from the perspective of signal detection theory and suggested that the dependency occurs because the observer's decision criterion systematically fluctuates throughout the experiment (Treisman and Williams, 1984; Dorfman and Biderman, 1971; Kac, 1962). More specifically, after a response, the observer shifts the criterion by a certain amount and the shift direction will determine the avoidance or repetition dependency. In the framework of drift-diffusion model, sequential dependency is suggested to happen because the starting point of evidence accumulation process is biased by the statistics of previous responses and/or feedback (Gold et al., 2008; Kim et al., 2017; de Lange et al., 2013). Along the same line, Akaishi et al. (2014) proposes that the history of previous decisions affects subsequent ones through a reinforcement learning process that updates based on subject's decision instead of feedback (note that there was no feedback in their experiments). In the context of Bayesian observer model, Angela and Cohen (2009) suggest that the observer assumes the environment is inherently volatile so that the probability of stimulus being repeated or alternated changes over time (see Glaze et al. (2015) for a similar account using drift diffusion model).

Essentially, the above explanations for sequential dependency posit that decisions across trials are dependent because they share a common source (response bias, decision criterion, history of previous trials/feedback or prior belief about environmental stability). That is different from what is typically thought about sequential dependency in the cognitive literature discussed above: a preceding decision directly biases subsequent decisions. The closest explanation in perceptual studies comes from Gold et al. (2008) and Akaishi et al. (2014) which suggest previous decisions directly affect subsequent ones. However, they also found that the effect of immediately preceding decision is either negligible or insufficient to account for the data. This discrepancy may be partly because in the perceptual studies, subjects only made a single judgment in each trial and the trials were often designed to be random and independent of each other. That differs from the cognitive literature in which subjects made multiple judgments in each trial. In the next section, I will review a small

but growing body of works in perception that bears a close resemblance to such situation.

1.2.2. Multiple-judgment task

In a perceptual task that is analogous to cognitive dissonance experiments, subjects had to make two sequential judgments on a cloud of moving dots: they first indicated whether the dots' motion direction was clockwise or counterclockwise of a decision boundary; then they had to report the perceived motion direction of the dots (Jazayeri and Movshon, 2007). The authors found that subjects' estimates in the second task were systematically biased away from the decision boundary. Importantly, the estimate bias is in the same direction as the decision. For instance, if their decision was clockwise in the first judgment, the estimate was biased away from the decision boundary in the clockwise direction. Furthermore, the bias magnitude was larger when the task was more difficult, that is, when the dots' motion direction was closer to the decision boundary or when the noise in the dots' motion direction was high. The findings were replicated in another study by Zamboni et al. (2016).

Along the same line, several studies on decision confidence found that decision may bias subsequent judgment of confidence. In a study by Zylberberg et al. (2012), subjects first made a categorical decision about the motion of a dots cloud (e.g. left or right), then they had to report their confidence on a continuous scale. The authors found that subjects' confidence judgments depend more on stimulus evidence supporting the decision than evidence opposing the decision. Similar results using face/house stimuli were found in Peters et al. (2017) at both behavioral and neural levels. The results imply that subjects are overconfident after the decision (see Yu et al. (2015a) for the under-confidence scenario).

Bronfman et al. (2015) show sequential dependency in a perceptual task similar to the sunk-cost fallacy scenario. Subjects were first presented with a sequence of 8 symbolic numbers and indicated whether the sequence average was less or greater than 50. Then another sequence of 8 numbers were displayed and subjects had to estimate the average of all 16 numbers. By regressing subjects' final estimates on the presented numbers, the authors found that subjects significantly down-weighted the evidence in the second sequence compared to the first sequence. They interpreted the results as showing that by making

a preliminary decision on the first sequence, subjects reduced sensitivity to the second sequence.

1.3. Bayesian decision theory

The key assumption of this theoretical approach is that humans maintain an internal model of how the world works (**generative/forward model**) and this internal model is characterized by probabilistic relations between variables. In this framework, perception can be formulated as an inverse process that infers hidden variables from the observed variables. For example, a baseball player wants to know the velocity of a ball flying in the air. In most cases, he doesn't have direct access to the true velocity of the ball (hidden variable) but instead he has some indirect information (observed variable) that is related to that variable. The most simple generative model assumes there is a true velocity v of the ball which the player doesn't know and would like to estimate (Fig. 1.1a). The player's visual system collects sensory evidence m_v about the true velocity which is a noisy copy of v . The player's task is to infer the ball's true velocity, v from the sensory evidence, m_v . Note that the internal models may vary in structure and level of details and that depends on several factors such as the task specification or subjects' past experience. In the baseball example, the player may have a more complicated internal model. For instance, his visual system may gather sensory evidence of the distance d traveled by the ball between two time points and the traveling time t (Fig. 1.1b). In this case, the player has to infer the ball's velocity from the sensory evidence of distance m_d and time m_t .

Among many ways to solve the problem, Bayesian decision theory provides a principled approach that can be interpreted intuitively. There are 3 main components of Bayesian decision theory: the prior distribution, the likelihood function and the loss function, each of which will be explained in detail below.

The **prior distribution** represents the observer's prior belief about the variable to be estimated before the sensory evidence is gathered. Technically, it is a probability distribution over the feasible space of the estimated variable. In the baseball example, the prior distribution $p(v)$ represents the probability of possible velocity values of the ball before the

player observes the current flying ball. An example of prior distribution shown in Fig. 1.1c illustrates the observer’s prior belief that the ball’s velocity tends to be slow rather than fast. Several studies provide supportive evidence for this tendency of people’s prior belief about speed in general (Weiss et al., 2002; Stocker and Simoncelli, 2006; Hürlimann et al., 2002).

Another important factor in Bayesian decision theory is the observer’s belief about the process that generates the sensory information. In the baseball example, this process may include how light is reflected off the ball and arrives at the player’s retina and how the retina and early visual system convert that incoming light pattern into neural signals. The whole process is often simplified by a noise model $p(m_v|v)$ that indicates the probability of obtaining a sensory evidence m_v given the true velocity v . In the estimation task, the sensory evidence is observed and fixed. When we write $p(m_v|v)$ as a function of v , it becomes the **likelihood function**. Therefore, the likelihood function represents the observer’s uncertainty about the sensory evidence and is totally determined by the whole process that takes place from the stimulus to the sensory evidence. For example, if the observer assumes the sensory evidence v_m is the result of adding Gaussian noise to the true velocity v , the conditional distribution $p(m_v|v)$ is a Gaussian distribution centered on v and the standard deviation is constant for all v . That gives rise to a likelihood function that has the Gaussian shape as illustrated in Fig. 1.1c.

Perception not only serves as a subjective experience but also has real consequences when it guides the observer’s action in most practical situations (e.g. whether a player successfully catches a ball or not). The observer’s belief about the consequence of its decision is quantified by the **loss function** which indicates how the decision error is penalized. Specifically, the loss function $L(\hat{v}, v)$ represents the cost of making decision \hat{v} when the true stimulus value is v . For computational convenience, the loss function is often defined as a parametric function of the difference between the decision and the true stimulus $\hat{v} - v$ (i.e. the decision error). Fig. 1.1c illustrates the squared loss function $L(\hat{v}, v) = (\hat{v} - v)^2$ which penalizes large error more strongly than small error. Another common loss function

is an absolute function of the error, $L(\hat{v}, v) = |\hat{v} - v|$ which also penalizes large error more but to a lesser extent than the squared loss. The loss function typically depends on the task and the reward structure. For instance, the squared or absolute loss functions may be appropriate for the baseball example because a small error is likely to be less costly than a large error (the player can still catch the ball if his hand is just a little bit away from the true position of the ball). In contrast, when the observer has to make a categorical decision (e.g. answering a multiple-choice question in an exam), the 0/1 loss function is often more suitable because getting credit requires an exactly correct answer:

$$L(\hat{v} - v) = \begin{cases} 0 & \text{if } \hat{v} = v \\ 1 & \text{if } \hat{v} \neq v \end{cases} \quad (1.1)$$

Given that all three main components (prior distribution, likelihood function and loss function) are well-defined in a task, Bayesian decision theory provides a **Bayes optimal rule** to make a decision. By using Bayes rule, the prior distribution is combined with the likelihood function to form the posterior distribution:

$$p(v|m_v) = \frac{p(m_v|v) \cdot p(v)}{p(m_v)} \quad (1.2)$$

In the baseball example, the posterior distribution $p(v|m_v)$ represents the player's belief about the ball's velocity v after integrating the sensory evidence with his prior belief. Because $p(m_v)$ is independent of v in Eqs. (1.2), it is just a normalizing constant for the purpose of estimating v . Then depending on the loss function, the observer makes a point estimate \hat{v} based on the posterior distribution $p(v|m_v)$. For instance, if we use the squared loss function, the estimate is the mean of the posterior distribution:

$$\hat{v}(m_v) = \int v \cdot p(v|m_v) dv \quad (1.3)$$

For the absolute and the 0/1 loss function, the estimates would be the median and the

mode of the posterior distribution, respectively. Note that the estimate \hat{v} is a function of the sensory evidence m_v . So Bayes optimal decision rule is a deterministic mapping between the sensory evidence and the estimate.

So far Bayesian decision theory has been described in terms of the observer's beliefs about the world (i.e. the generative model, prior, likelihood and loss function) and how the observer makes Bayes optimal decision with these beliefs. An obvious question is whether those beliefs reflect the facts about the physical world. It is an interesting question because if the observer's beliefs match the physical world, Bayesian decision theory is not just a **descriptive framework** which only aims to describe what is going on but it is also a **prescriptive framework** which defines the optimal strategy to make a decision. In fact, this issue has been investigated long before the Bayesian approach to perception becomes popular. Brunswik and Kamiya (1953) studied whether Gestalt rules of how humans group items accord with the statistical regularities of the environment. They found that some rules indeed reflect the statistical regularities of the world (e.g. the proximity principle in which items closer to each other tend to belong to the same object). Several studies have followed up on that using Bayesian approach and also found a fairly good match between human beliefs and the statistical structure of the environment (Geisler et al., 2001; Geisler and Perry, 2009; Burge et al., 2010; Girshick et al., 2011; Kim and Burge, 2018). Those studies mostly touch on the prior distribution. Other components of Bayesian decision theory (generative model, likelihood and loss function) seem to be harder to investigate. For example, it is not always straightforward to determine the loss function. Gambles in the casino (roulette, blackjack, etc.) may have a clear reward structure. In contrast, choosing between orange and apple primarily depends on individual preferences. As a result, there is still a lot of open questions regarding Bayesian decision theory as a prescriptive framework.

Application of Bayesian decision theory to the study of perception has been both successful and fruitful because it can account for a wide range of experimental data and has the potential to generate interesting testable predictions. Several studies have shown that human perceptual decisions agree to a large extent with the predictions of Bayesian inference

across a wide variety of stimuli and tasks such as the perception of depth (Ernst and Banks, 2002; Knill, 2007; Jacobs, 1999), velocity (Ascher and Grzywacz, 2000; Hürlimann et al., 2002), color (Brainard and Freeman, 1997; Brainard et al., 2006), sound (Gifford et al., 2014), time (Jazayeri and Shadlen, 2010; Acerbi et al., 2012; Shi et al., 2013), multi-sensory stimulus (Körding et al., 2007) as well as sensorimotor tasks (Körding and Wolpert, 2004; Körding and Wolpert, 2006).

1.4. Supplemental-psychophysics methods

Scaling procedure: In a common form of scaling, subjects are typically asked to assign a numerical value to their perceived intensity of the stimulus (Stevens, 1946, 1957; Marks, 1974). For example, the experimenter presents a standard sound as an anchor and assign a number 10 to the loudness of that sound. Then another comparison sound is played and subjects have to indicate how loud the comparison is relative to the standard sound (e.g. 20 if the comparison is twice as loud as the standard) (Stevens, 1956). Not all scaling methods involve numerical assignment. In the Farnsworth-Munsell 100-Hue test, subjects are shown a row of colored papers that are randomly ordered and vary in hue. The first and last papers are fixed and subjects have to order the other papers to make a "smooth color series" (Farnsworth, 1943).

Forced choice task: Because many scaling methods involve subjects' comprehension or production of numbers, any abnormality in number perception may contaminate the result. The forced choice method may overcome this by asking subjects to make a simple categorical response to the stimulus (Fechner, 1860; Green and Swets, 1966). In a typical two-alternative forced choice experiment, two stimuli (e.g. two circular gratings with different orientations) are presented and subjects have to make a binary decision (e.g. which grating is more clockwise). By analyzing the data of forced choice method using signal detection theory (Macmillan and Creelman, 2004), we can make useful inference about subjects' sensitivity to the stimulus (e.g. subjects' discrimination threshold) and subjects' decision biases (e.g. subjects' tendency to respond clockwise regardless of the stimulus).

Adjustment/Matching method: A classic experiment using this method was performed more

than a hundred years ago by James Clerk Maxwell (Maxwell, 1857). To study human color perception, he constructed a countertop which consists of two disks put on top of each other. The smaller disk on top is painted with a reference color. The bigger disk on the bottom is divided into three wedges whose colors are vermilion, emerald green and ultramarine blue (the three primary colors). During the experiment, the disks are spun very fast so that the three colors of the big disk on the bottom are seen as being blended together, resulting in the percept of a single color. Subjects were tasked to adjust the proportion of three primary colors until the color of the top and the bottom disks are the same. An illustration of the design can be found here (<http://www.handprint.com/HP/WCL/colortop.html>).

There are many variations of the above methods and interested readers can find more details in Kingdom and Prins (2010), Macmillan and Creelman (2004) and Lu and Doshier (2013).

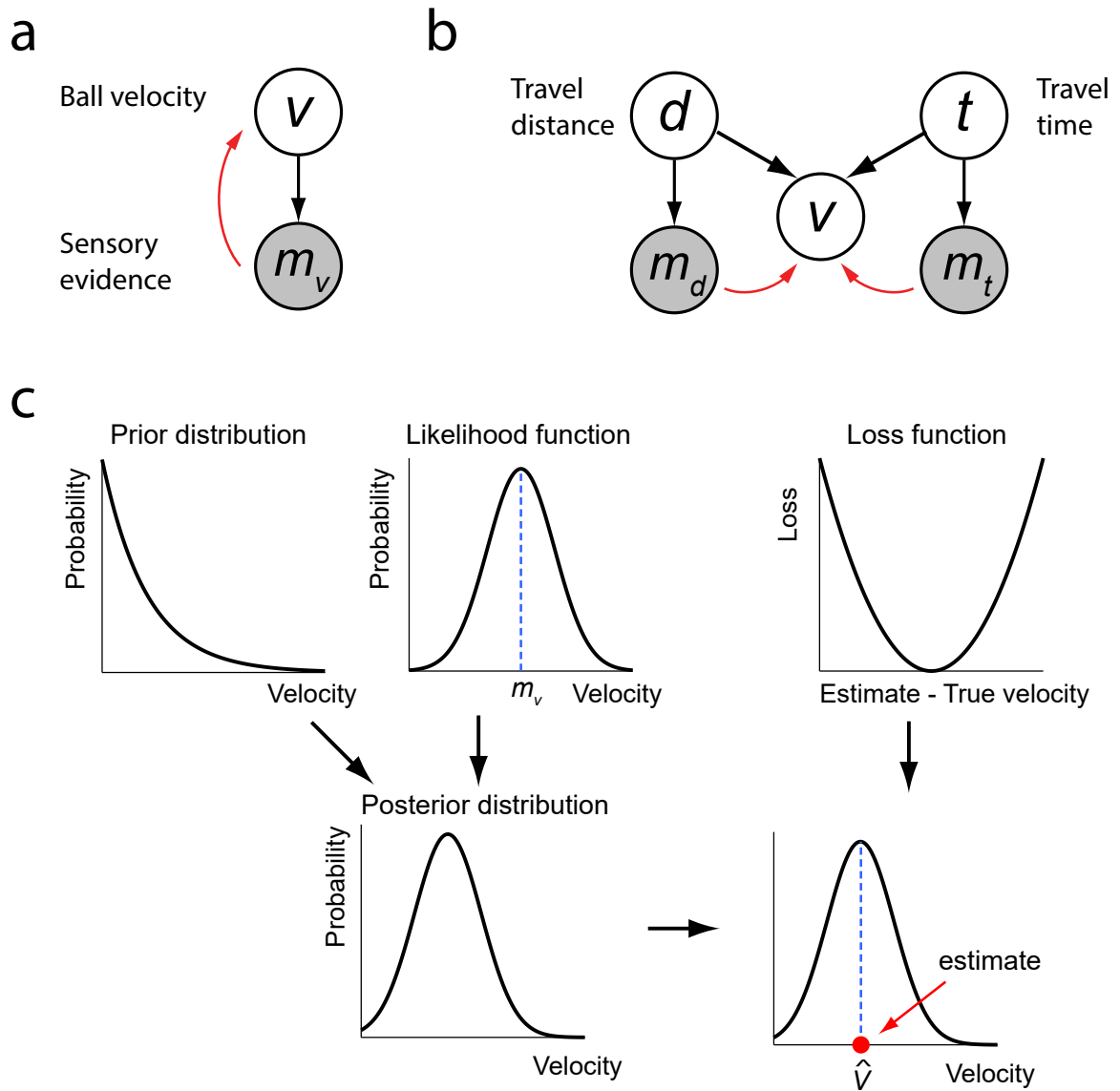


Figure 1.1: *Example to illustrate Bayesian decision theory.* (a) A simple generative model can be used to represent the problem of estimating the velocity of a ball: v indicates the true velocity and m_v indicates the sensory evidence about the ball's velocity. The task is to infer v from m_v . (b) A more complicated generative model can be used. Here the observer infers the ball velocity from sensory information about the distance traveled by the ball and the travel time. (c) Three main components of Bayesian decision theory are shown: prior distribution, likelihood function and loss function. The prior distribution is an exponentially decreasing function, indicating the observer's belief of predominantly slow moving objects. The likelihood function is a Gaussian distribution centered on the sensory measurement. Loss function is a squared function of the difference between estimated and true velocity. A Bayesian observer computes the posterior distribution from the prior distribution and the likelihood function. Then an estimate is made by computing the mean of the posterior distribution.

Chapter 2

Self-consistent inference

2.1. Introduction

In this chapter, I set the theoretical foundation for the thesis by refining and experimentally validating a hypothesis that may explain for several sequential dependencies between decisions discussed in Chapter 1 (see section 1.2 and 1.3.2). Specifically, I tested the hypothesis that an observer’s inclination to maintain a self-consistent, hierarchical interpretation of the world leads to the observed post-decision biases in perceptual judgments. We expressed this hypothesis with a *self-consistent Bayesian observer model*, which assumes that a subject’s estimate is not only conditioned on the sensory evidence but also on the subject’s preceding categorical judgment. The model extends a previous formulation (Stocker and Simoncelli, 2007) as it takes into account that sensory information in working memory degrades over time, which has important implications with regard to the behavioral benefits of the proposed hypothesis. I validated the model with three different psychophysical experiments that were based on a perceptual task sequence in which subjects were asked to perform a categorical orientation judgment followed by an orientation estimate of a visual stimulus. Results from the first experiment, using a similar design as the original experiment by Jazayeri and Movshon (2007), suggest that post-decision biases are general to different low-level visual stimuli. The other two experiments were targeted to specifically test key assumptions of the model such as the probabilistic dependence of the biases on sensory priors. I show that the self-consistent observer model accurately accounts not only for the data from our experiments but also from a previous study (Zamboni et al., 2016). The similarity to well-documented bias effects in social psychology and economics suggests that the proposed self-consistent model may provide a general framework for modeling sequential

decision-making¹.

2.2. Replicating choice-induced bias

The first goal was to obtain a quantitative characterization of how the perceptual inference process is affected by a preceding, categorical decision task (Fig. 2.1a). In Experiment 1 subjects first indicated whether the overall orientation of an array of lines was clockwise (cw) or counter-clockwise (ccw) relative to a discrimination boundary (discrimination task), and then had to reproduce their perceived overall stimulus orientation by adjusting a reference line (estimation task). The experimental design was similar to that of a previous experiment (Jazayeri and Movshon, 2007) (see also Zamboni et al. (2016)) with the notable exceptions that we used an orientation rather than a motion stimulus, and that subjects had to perform the estimation task in every trial rather than only in one third of the trials (Fig. 2.1b). We found that subjects' perceptual behavior was very similar to the findings of these previous studies (Jazayeri and Movshon, 2007; Zamboni et al., 2016). Discrimination performance monotonically depended on stimulus noise. Furthermore, reported stimulus orientations showed clear repulsive biases away from the discrimination boundary with larger biases for larger stimulus noise and orientations closer to the boundary (Fig. 2.1c). This behavior is fully revealed in the overall distributions of subjects' estimates for every noise level (Fig. 2.1d; see Figure supplement 2.11 for individual subjects): The distributions exhibit a characteristic bimodal shape for stimulus orientations close to the decision boundary, with subjects' estimates biased away from the decision boundary toward the side that corresponds to their preceding discrimination judgment. Together with previous findings (Jazayeri and Movshon, 2007; Zamboni et al., 2016), the results of Experiment 1 suggest that the observed post-decision biases are independent of the specific type of stimulus used. Also, they indicate that subjects' anticipation of the frequency of the estimation task (every trial *vs.* only in 1/3 of the trials) does not play a role in causing the biases.

¹This is in line with the general idea by Lecky (1945) that an individual must remain a self-consistent unit in the way it interacts with the outside world.

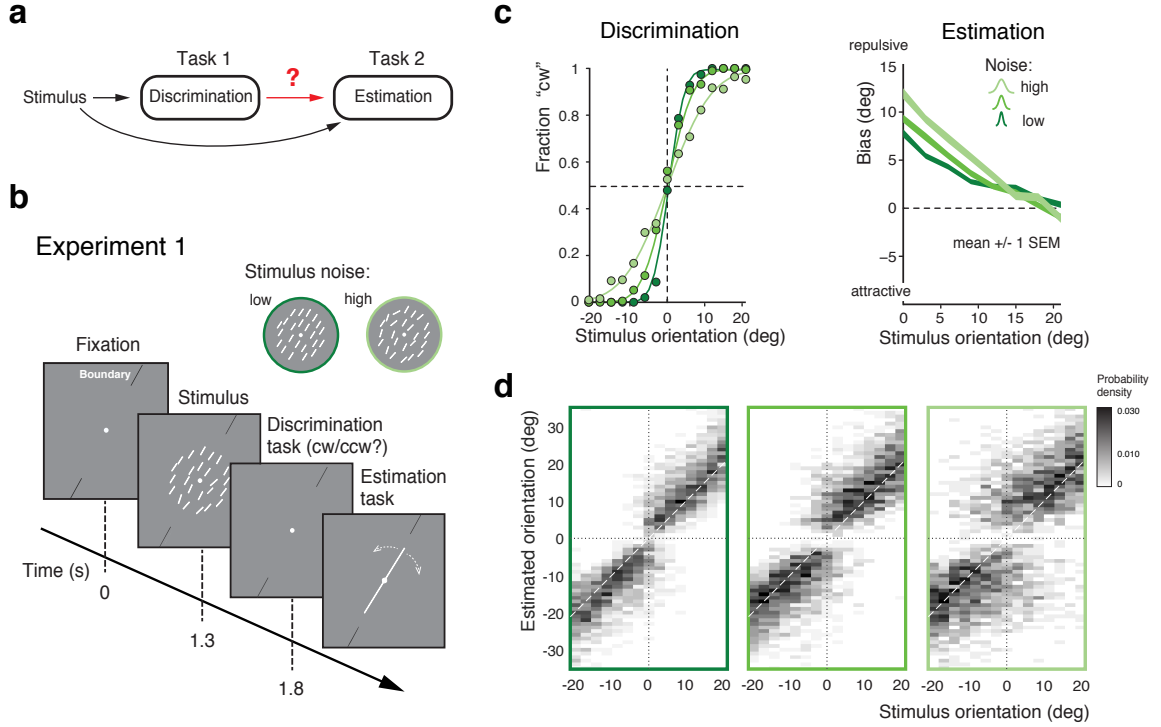


Figure 2.1: *Post-decision biases in a perceptual task sequence.* (a) Perceptual decision-making in an estimation-discrimination task sequence: Does a discrimination judgment causally affect a subject's subsequent perceptual estimate? (b) Experiment 1: After being presented with an orientation stimulus (array of lines), subjects first decided whether the overall array orientation was clockwise (cw) or counter-clockwise (ccw) of a discrimination boundary, and then had to estimate the actual orientation by adjusting a reference line with a joystick. Different stimulus noise levels were established by changing the orientation variance in the array stimulus. (c) Psychometric functions and estimation biases (combined subject). Estimation biases are only shown for correct trials and are combined across cw and ccw directions. Subjects show larger repulsive biases the noisier the stimulus and the closer the stimulus orientation was to the boundary. (d) Distributions of estimates for the three stimulus noise levels tested, plotted as a function of the true stimulus orientation relative to the discrimination boundary (combined subject). Estimates are clearly biased away from the discrimination boundary forming a characteristic bimodal pattern.

2.3. Self-consistent Bayesian observer model

How are these post-decision biases explained? A Bayesian observer that regards the task sequence as a two independent inference processes does not predict the biases. The observer uses Bayesian statistics to determine the correct categorical judgment (*e.g.* 'cw') based on the stimulus response of a population of sensory neurons, and does the same to infer the best possible estimate of the stimulus orientation (Fig. 2.2a). Consequently, this observer's discrimination judgment does not affect the estimation process; orientation estimates are unimodally distributed around the true stimulus orientation and do not exhibit the characteristic bimodal pattern that we have observed in Experiment 1 (Fig. 2.1d). In the context of this paper, we refer to this model as the "independent" observer.

In contrast, we propose an observer model that regards the two tasks as causally dependent (Fig. 2.2b); that is, by making the discrimination judgment the observer constrains its subsequent estimation process to consider only those stimulus orientations that are consistent with the judgment. It is as if the observer regards its own, subjective discrimination judgment as an objective fact. Such behavior seems irrational (obviously, the judgment could be incorrect) and furthermore leads to characteristic estimation biases away from the discrimination boundary. It has the advantage, however, that the observer's perceptual inference process remains self-consistent throughout the entire task sequence at any moment in time. We refer to this model as the "self-consistent" observer. It can be formulated as a conditioned Bayesian model (Stocker and Simoncelli, 2007) that jointly accounts for subjects' behavior in both the discrimination and the estimation task. However, a closer comparison between the predicted (Fig. 2.2b) and the measured distribution of the estimates (Fig. 2.1d) reveals that this basic formulation does not capture all details accurately.

We formulated the self-consistent observer model as a two-step inference process over the extended hierarchical generative model shown in Fig. 2.3a: Based on a noisy sensory signal m , the observer first infers the category C ('cw' or 'ccw') by performing the discrimination task and then infers the stimulus orientation θ in the estimation task. Because the stimulus has long disappeared by the time the observer performs the estimation task, we

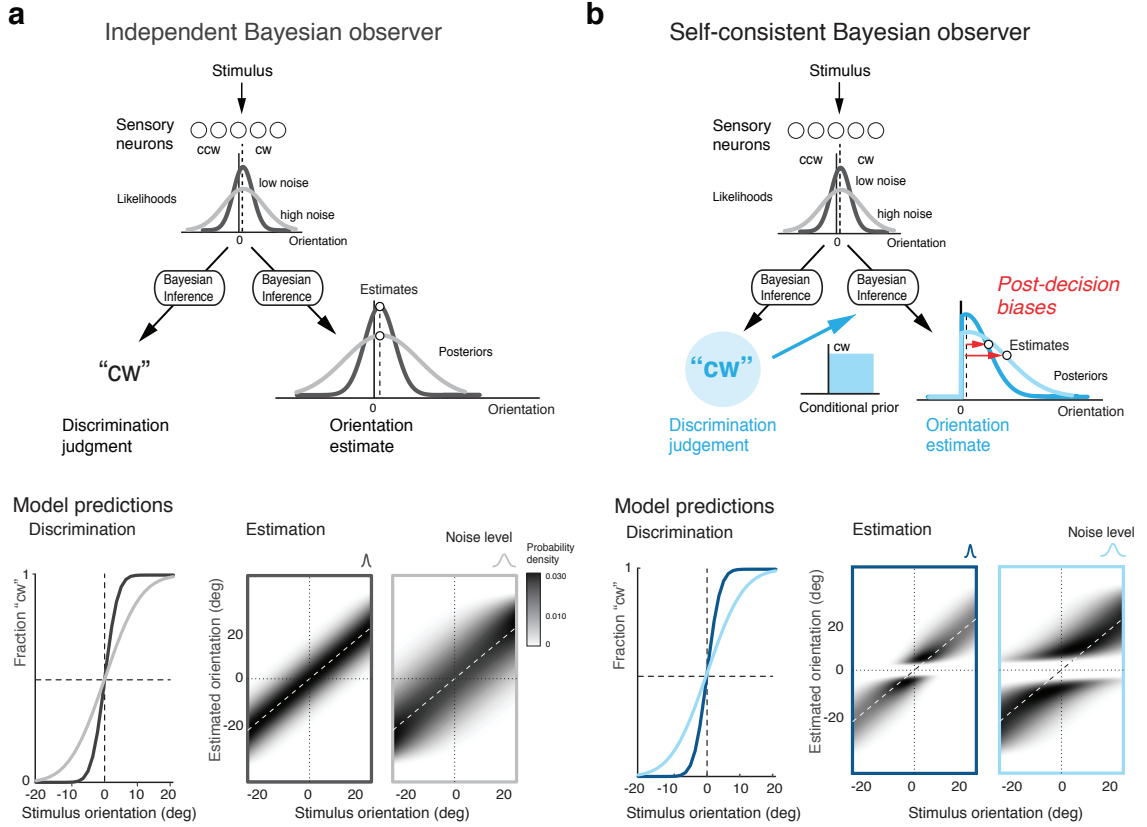


Figure 2.2: *Bayesian observer models for the perceptual task sequence.* (a) The discrimination judgment does not affect the estimated stimulus orientation for an observer who considers both tasks independently. (b) In contrast, the self-consistent observer imposes a causal dependency such that the judgment in the discrimination task (e.g. 'cw') conditions the estimation process in form of a choice-dependent prior. It effectively sets the posterior probability to zero for any orientation value that is inconsistent with the preceding discrimination judgment. The truncated posterior distribution, together with a loss function that penalizes larger estimation errors stronger than smaller ones, leads to the characteristic bimodal distribution pattern. Note, however, that this basic formulation is not quite sufficient to explain some details of the distribution pattern of subjects' estimates (Fig. 2.1d).

assume that estimation of θ must rely on a noisy memory recall m_m of the sensory signal m . Inference on θ is then conditioned on the preceding discrimination judgment (*e.g.* $C = \text{'cw'}$), which ultimately results in the characteristic repulsive estimation biases. Finally, we also took into account that subjects' report of their perceived stimulus orientation is corrupted by motor noise. We measured motor noise for every subject in a control experiment (see Figure supplement 2.12) and subsequently used these measured values for all model fits and comparisons. The self-consistent observer model provides a full account of both the observer's discrimination judgment and orientation estimate in each trial and is thus jointly predicting a subject's psychometric function as well as the distribution of their orientation estimates.

Figure 2.3b shows the model fit to the data from Experiment 1 for the combined subject. The stimulus noise level determines both the slope of the psychometric curves in the discrimination task and the magnitude of the bias in the estimation task, which is well predicted by the model. A comparison between the distributions of subjects' estimates and model estimates fully reveals the extent to which the model accurately accounts for the observed human behavior (Fig. 2.3c; See Figure supplement 2.13 for a histogram representation). Note that all model predictions in this paper are the result of a joint fit to both the measured psychometric functions of the discrimination task and the estimation distribution.

The self-consistent observer model can also account for the substantial individual differences in behavior across subjects. While individual bias patterns are all repulsive, they vary across subjects both in shape and magnitude (Fig. 2.4a). These variations are well captured by the model and reflected in individual differences in the fit parameter values such as the prior width and the level of sensory noise (Fig. 2.4b; see Figure supplement 2.14 for a goodness-of-fit analysis). Interestingly, all subjects seemed to substantially overestimate the width of the stimulus prior compared to the true stimulus distribution. This did not come entirely as a surprise because subjects were never explicitly informed about the stimulus range and thus had to learn it over the course of the experiment. Consistent with this interpretation is the fact that the subject with the most accurate estimate of the

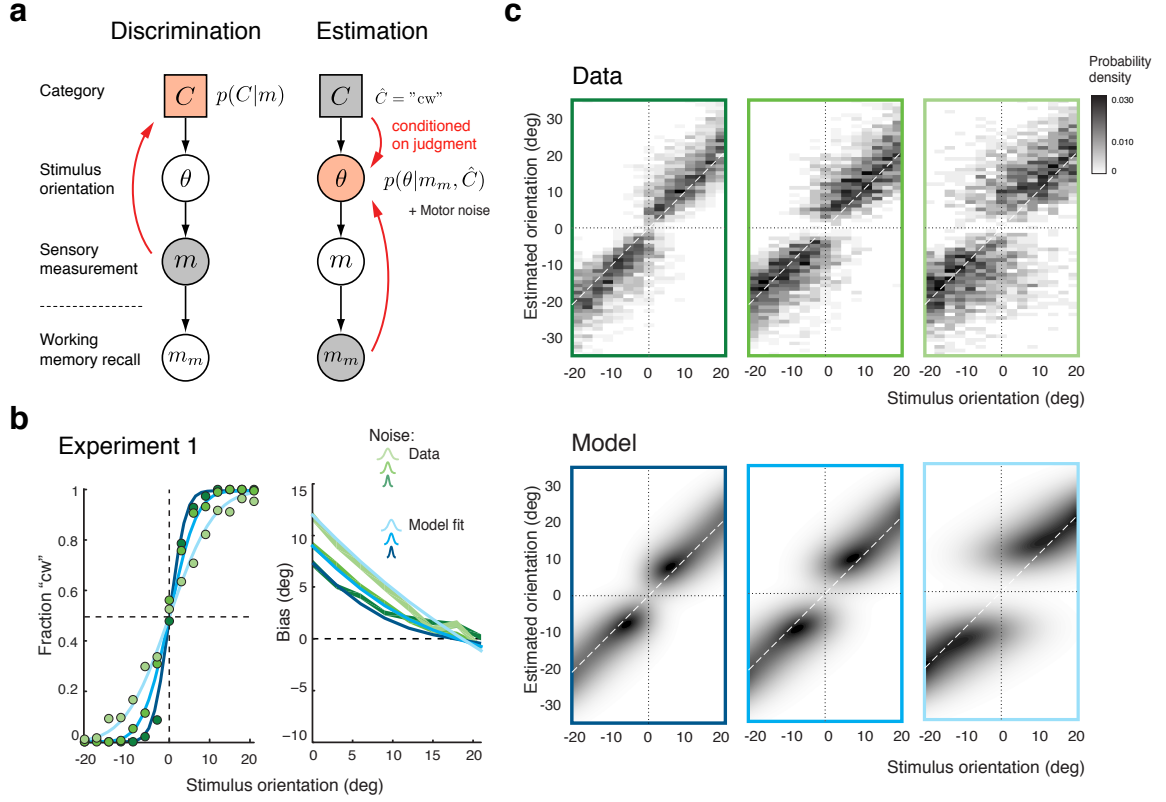


Figure 2.3: *The self-consistent Bayesian observer model.* (a) Directed graph representing the generative hierarchical model: Sensory measurement m is a noisy sample of stimulus orientation θ . Every θ belongs to one of two categories $C = \{\text{'cw'}, \text{'ccw'}\}$. Given an observed m , the self-consistent model first performs inference over C (discrimination task), and then infers the value of θ *conditioned* on the preceding discrimination judgment (*e.g.* $\hat{C} = \text{'cw'}$) (estimation task). Inference for the estimation task is assumed to be based on a noisy memory recall m_m of the sensory measurement m . Conditioning on the categorical choice sets the posterior $p(\theta|m_m, \hat{C})$ to zero for all values of θ that do not agree with the choice. This shifts the posterior probability mass away from the discrimination boundary and results in the repulsive post-decision biases for any loss function that more strongly penalizes large errors than small ones. Because subjects were instructed to provide estimates as accurate as possible we assumed a loss function that minimizes mean squared error (L_2 loss). (b) We jointly fit the observer model to all discrimination-estimation data pairs of the combined data across all subjects in Experiment 1 (combined subject). (c) The model not only predicts the mean estimation bias (as shown in (b)) but also the entire distributions of estimates, including those trials where discrimination judgments were incorrect. Data and model show the characteristic bimodal pattern for orientation estimates. Each column corresponds to one of the three stimulus noise conditions.

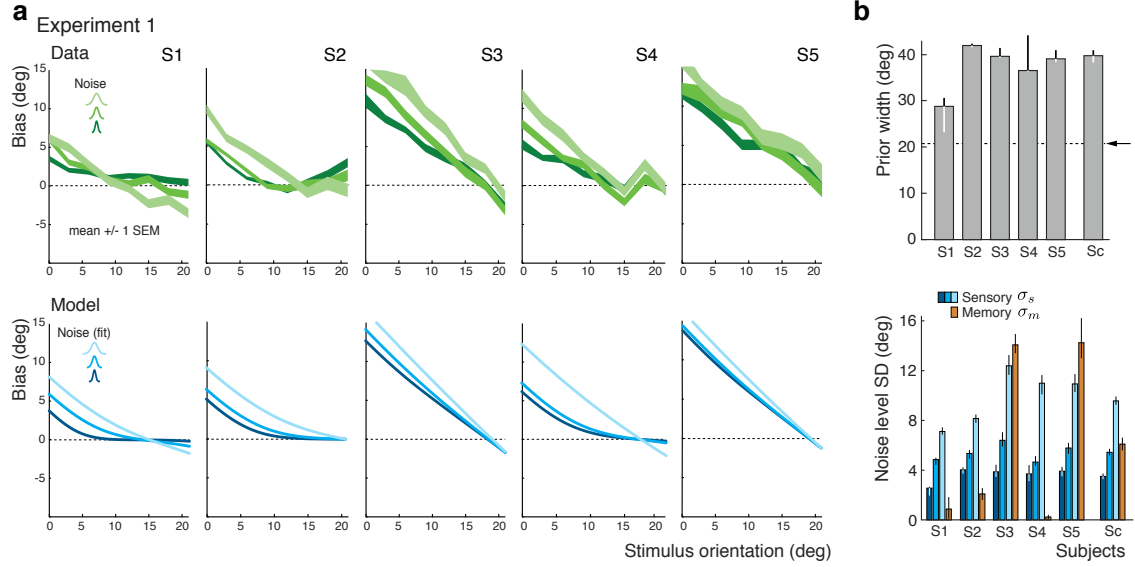


Figure 2.4: *Experiment 1: Data and model fits for individual subjects.* (a) Individual subjects (S1 non-naïve) showed substantial variations in their bias patterns (green curves). These variations are well explained by individual differences in the fit parameter values of the self-consistent model (blue curves). For example, the width of the prior directly determines the location where the bias curves intersect with the x-axis. (b) Fit prior widths w_p and noise levels for the five individual subjects plus the combined subject (Sc). Subjects' prior widths suggest that they consistently overestimated the actual stimulus range in the experiment (± 21 degrees; arrow). For all subjects, fit sensory noise σ_s was comparable and monotonically dependent on the actual stimulus noise. Memory noise σ_m was mostly small as expected, yet dominated for subjects S3 and S5. These two subjects performed poorly in the estimation task, suggesting that they were not trying to provide an accurate orientation estimate but simply pointed the cursor to roughly the middle of the stimulus range on the side of the discrimination boundary they picked in the discrimination task. The resulting bias curves are basically independent of the stimulus noise and have a slope of approximately -1. The model captured this behavior by assuming that the sensory information was “washed out” with a large amount of memory noise. The full model also contained a motor noise component that was determined for each subject in a separate control experiment. All errorbars represent the 95% confidence interval computed over 100 bootstrapped sample sets of the data. See Methods section 2.8 for details.

prior distribution was the only non-naïve subject S1 who had plenty of extra exposure to the stimulus range through the participation in various pilot experiments. Extracted noise levels differ across subjects with the worst subject being approximately twice as noisy as the best (non-naïve subject S1), yet consistently increase for increasing stimulus noise levels.

2.4. Validating self-consistent Bayesian observer model

2.4.1. Key features of the model

We ran two additional experiments that were designed to specifically probe two key features of the self-consistent model: Experiment 2 was aimed at testing how subject’s orientation estimates were dependent on their precise knowledge of the stimulus prior and thus were consistent with the results of Bayesian inference; Experiment 3 examined whether subjects indeed treated their discrimination judgments as if they were correct. We recruited a new set of subjects (S6-S9, plus S1) that performed both experiments. By jointly fitting the data from both experiment, we also tested how well the model can generalize across tasks².

Experiment 2 was identical to Experiment 1 except that at the beginning of each trial, subjects were explicitly reminded of the total range within which the stimulus orientation would occur in the trial (Fig. 2.5a). Our assumption was that an explicit display of the stimulus range provided subjects with a better and presumably narrower representation of the stimulus distribution (given that subjects seemed to substantially overestimate the prior in Experiment 1). If so, then the self-consistent observer model would predict a shift of the bias curves’ crossover point towards the discrimination boundary. As shown in Figs. 2.5b, the measured bias curves indeed show the predicted shift compared to the bias curves measured in Experiment 1 (Fig. 2.3b). This shift is also clearly visible in the distributions of the orientation estimates (see Figure supplement 2.15 for distributions of individual subjects), which is again accurately accounted for by the model (Fig. 2.5c).

In Experiment 3, we separated the discrimination judgment from the discrimination task. Subjects were no longer asked to perform the orientation discrimination task but

²Because subject S1 (the only non-naïve subject) performed all three experiments, fits and model parameters for this subject are the result of a joint fit to the data from all three experiments.

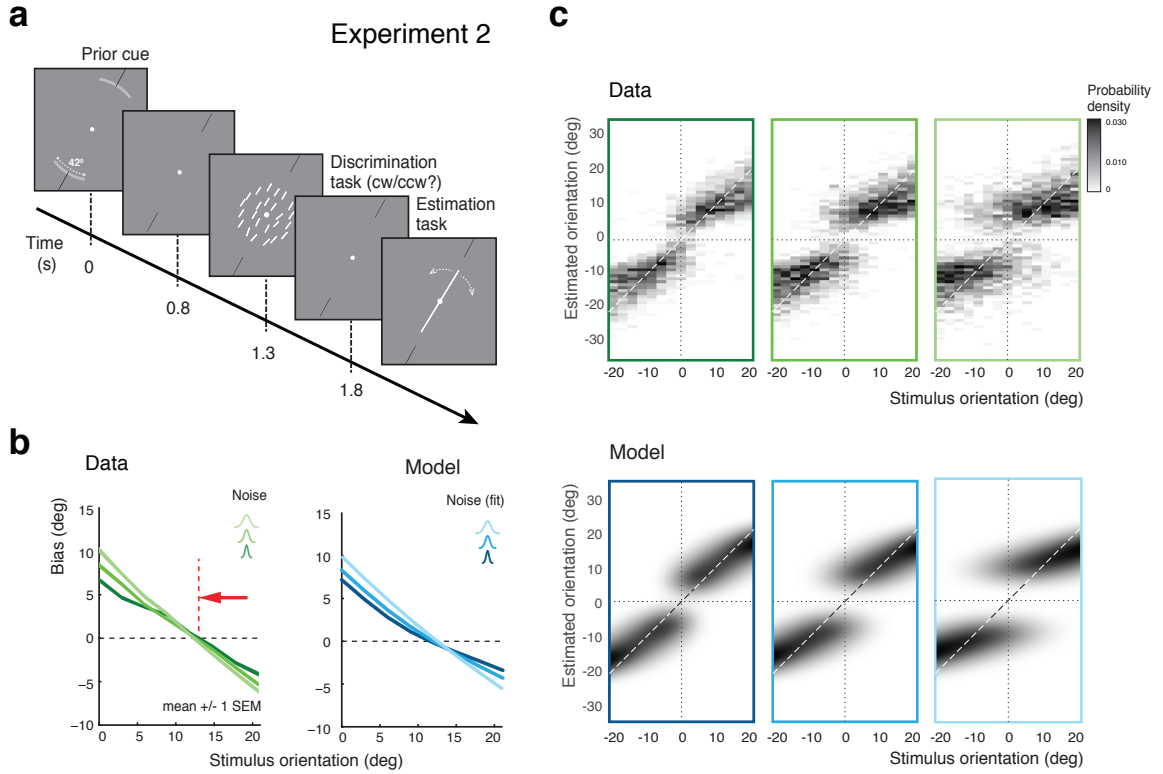


Figure 2.5: *Effect of the stimulus prior*. (a) Experiment 2 was identical to Experiment 1 except that at the beginning of each trial, subjects were shown the total range within which the stimulus orientation would occur in the trial (gray arc, subtending ± 21 degrees). (b) We hypothesize that reminding subjects of the exact stimulus range at the beginning of each trial helps them to form a more accurate (and more narrow) representation of their stimulus prior. If subjects' orientation estimates were indeed the result of the conditioned Bayesian inference as assumed by the self-consistent observer model, then the bias curves should shift towards the discrimination boundary. The data support this prediction: Subjects' bias curves (combined subject, see Fig. 2.7 for individual subjects) are shifted towards the discrimination boundary compared to Exp. 1. (c) As with Exp.1, the fit self-consistent model provides an accurate description of the distribution pattern of subjects' orientation estimates.

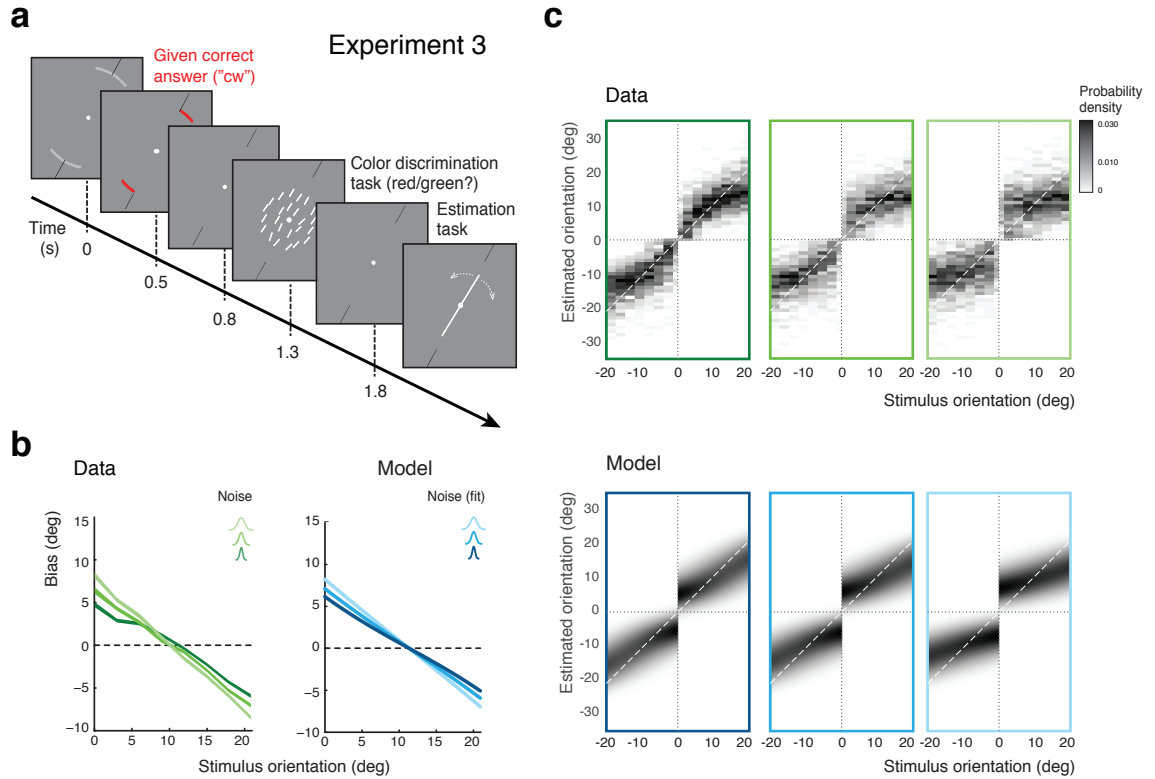


Figure 2.6: *Self-made vs. given category assignment.* (a) Experiment 3: Instead of performing the discrimination judgment themselves, subjects were provided with a cue indicating the correct category assignment right before the stimulus was presented. Then, after stimulus presentation, subjects first performed an unrelated color discrimination task in place of the orientation discrimination task (they needed to remember the randomly assigned color (red/green) of the cue indicating the correct category) before indicating their perceived stimulus orientation. (b) According to our model we should see similar estimation biases in Exp. 2 and 3, which is indeed what we found: because the data are from the same subjects we can directly compare the bias curves with the results from Exp. 2. (c) Again, the fit model well accounts for the overall distribution of orientation estimates (combined subject; see Figure supplement 2.16 for distributions for individual subjects)). Because the discrimination judgment was given and always correct independent of the noise in the sensory measurement m , estimates only occurred in the “correct” quadrants. For the same reason the model also predicts slightly smaller bias magnitudes (compared to Exp.2), which is also matched by the data (see also Fig. 2.7b).

instead were signaled right at the beginning of each trial whether the stimulus orientation would be cw or ccw (Fig. 2.6a). Subjects were instructed that this categorical information was always correct, which it was. They then performed an unrelated color discrimination task before finally performing the estimation task. The self-consistent model predicts estimation biases that are basically identical to those of Experiment 2 because it assumes that subjects treat their own judgment as correct when performing the estimation task. Indeed, as shown in Fig. 2.6b, subjects' estimation biases are very similar to the biases in Experiment 2 (Fig. 2.5b). Because both Experiment 2 and 3 were conducted on the same set of subjects, the results are directly comparable. We can rule out that subjects may have ignored the given category assignment in Experiment 3 and implicitly performed the orientation discrimination task instead. If this were the case, then subjects would have exhibited a large fraction of inconsistent trials (*i.e.* trials in which the estimated orientation was not in agreement with the given correct answer) in particular for orientations close to the discrimination boundary. This was not the case as we observed only small fractions of inconsistent trials (4% on average) that were of similar magnitude as the error rates for the (irrelevant) color discrimination task (2%). We discuss these inconsistent trials in more detail in the next section below.

We again extended our analysis to individual subjects' behavior. Figure 2.7a shows subjects' estimation biases in both experiments as well as the corresponding model predictions based on a joint fit to data from both Experiment 2 and 3. Bias patterns, while quite variable across subjects, are consistent across the two experiments for each subject. This confirms that the impact of the categorical discrimination judgment on the perceived orientation does not depend on whether the judgment was performed by the subjects themselves or not. The model captures both the inter-subject as well as the within-subject variability across the two experiments. Biases are slightly smaller in Experiment 3 compared to Experiment 2 for stimulus orientations close to the boundary. As shown in Fig. 2.7b, the model correctly predicts this difference because the self-made discrimination judgments in Experiment 2 are based on the noisy stimulus measurement m and therefore can be incorrect,

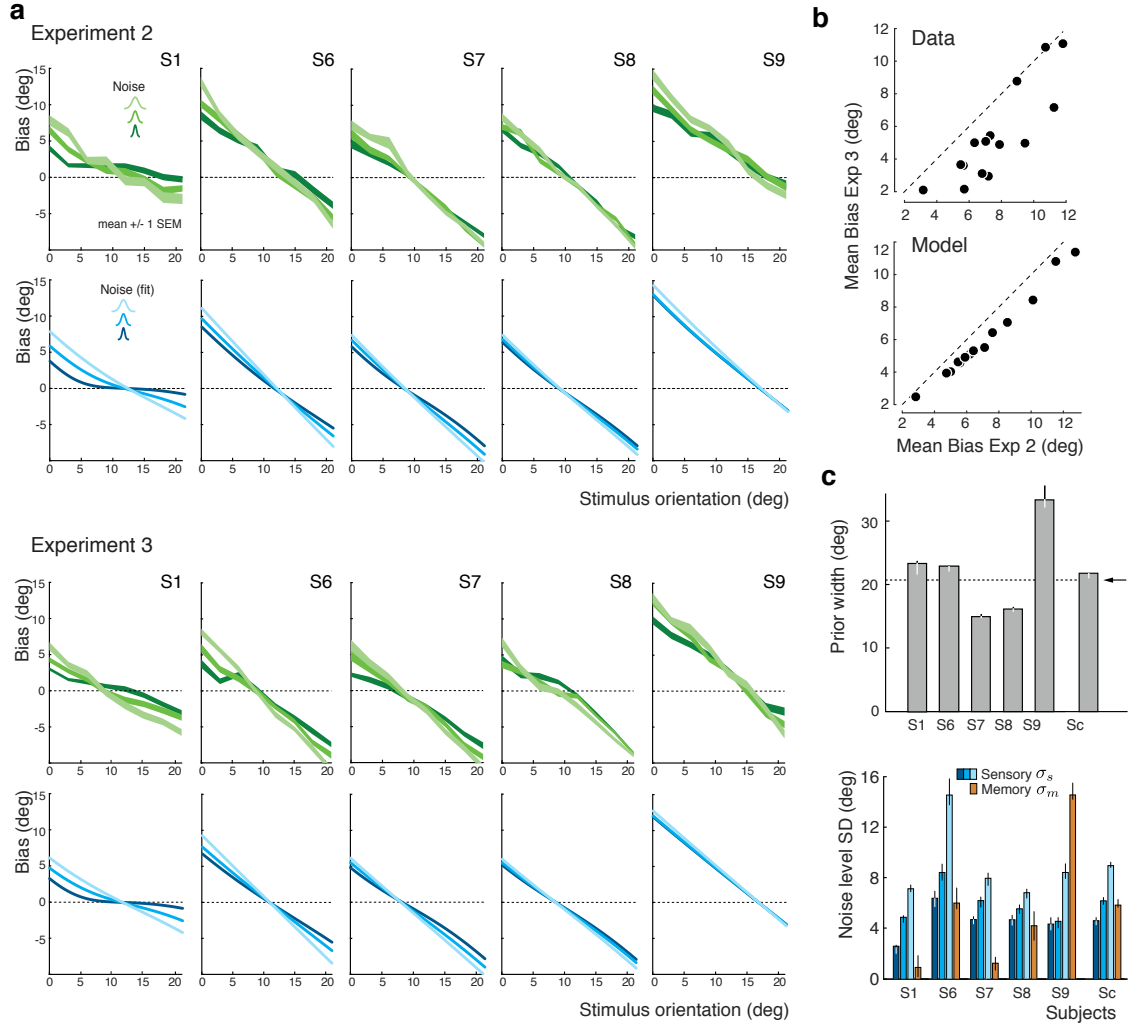


Figure 2.7: *Experiments 2 and 3: Joint fit to data for individual subjects.* (a) Five subjects (S1, S6-9) participated both in Exp. 2 and 3. We performed a joint model fit to the data from both experiments for every subject. Each column shows data (green curves) and model fit (blue curves) for a particular subject. As in Exp. 1, the bias pattern across subjects shows substantial variability yet is strikingly similar between the two experiments. (b) Comparing the mean biases observed in Exps. 2 and 3 reveals that biases in Exp. 3 are slightly smaller for stimulus orientations close to the boundary. This effect is predicted by the model. (c) Prior widths w_p and noise levels from the joint fit for individual subjects and the combined subject. Subjects' priors were closer to the experimental distribution than in Exp. 1 because in Exps. 2 and 3 subjects were reminded about the stimulus range at the beginning of each trial. Noise levels were comparable to those in Exp. 1 (for S1 we jointly fit data from all three experiments). Errorbars indicated the 95% confidence interval over 100 bootstrapped samples of the data. See Figure supplement 2.17 for a goodness-of-fit analysis.

while the category cues in Experiment 3 were always correct. Consequently, the predicted bias curves for Experiment 2 only represent trials for which the sensory measurement m was in favor of the correct judgment (*i.e.* m was on the correct side of the boundary) whereas the bias curves for Experiment 3 are computed over all trials. As a result, the biases in Experiment 3 are smaller for stimulus orientations for which there is a substantial chance that the noise pushes the measurements m to the other side of the discrimination boundary. Figure 2.7c shows the fit parameter values for individual subjects. Compared to Experiment 1, the subjective prior widths are substantially smaller and closer to the true stimulus prior width suggesting that explicitly reminding subjects of the true stimulus distribution in every trial was effective. As in Experiment 1, subjects showed large variations in subjective noise levels although they consistently were monotonic in actual stimulus noise. Fit memory noise levels were relatively small with the notable exception of subject S9 whose poor performance in the estimation task, quite similar to subjects S3 and S5 in Experiment 1, was picked up by the memory noise parameter.

2.4.2. Inconsistent trials

In a small fraction of trials (on average 4%) subjects' discrimination judgment and their subsequent orientation estimate were not consistent. We can show that these inconsistent trials are not a violation of self-consistent inference but rather can be entirely explained by two common sources of behavioral errors not related to perceptual inference: lapses and motor noise. In fact, we can accurately predict the estimation patterns and individual fractions of inconsistent trials based on model fits to the consistent data, and individual measurements of lapse rates and motor noise.

Figure 2.8a shows the distribution of subjects' estimates (combined data across all subjects and all stimulus noise conditions) in inconsistent trials for all three experiments, together with predictions from the self-consistent observer model for each of the two error sources. Lapses are defined as trials in which subjects by mistake pushed the wrong button in the discrimination judgment (respectively, incorrectly remembered the answer cue in Experiment 3) yet followed up with an orientation estimate that corresponded to their actual

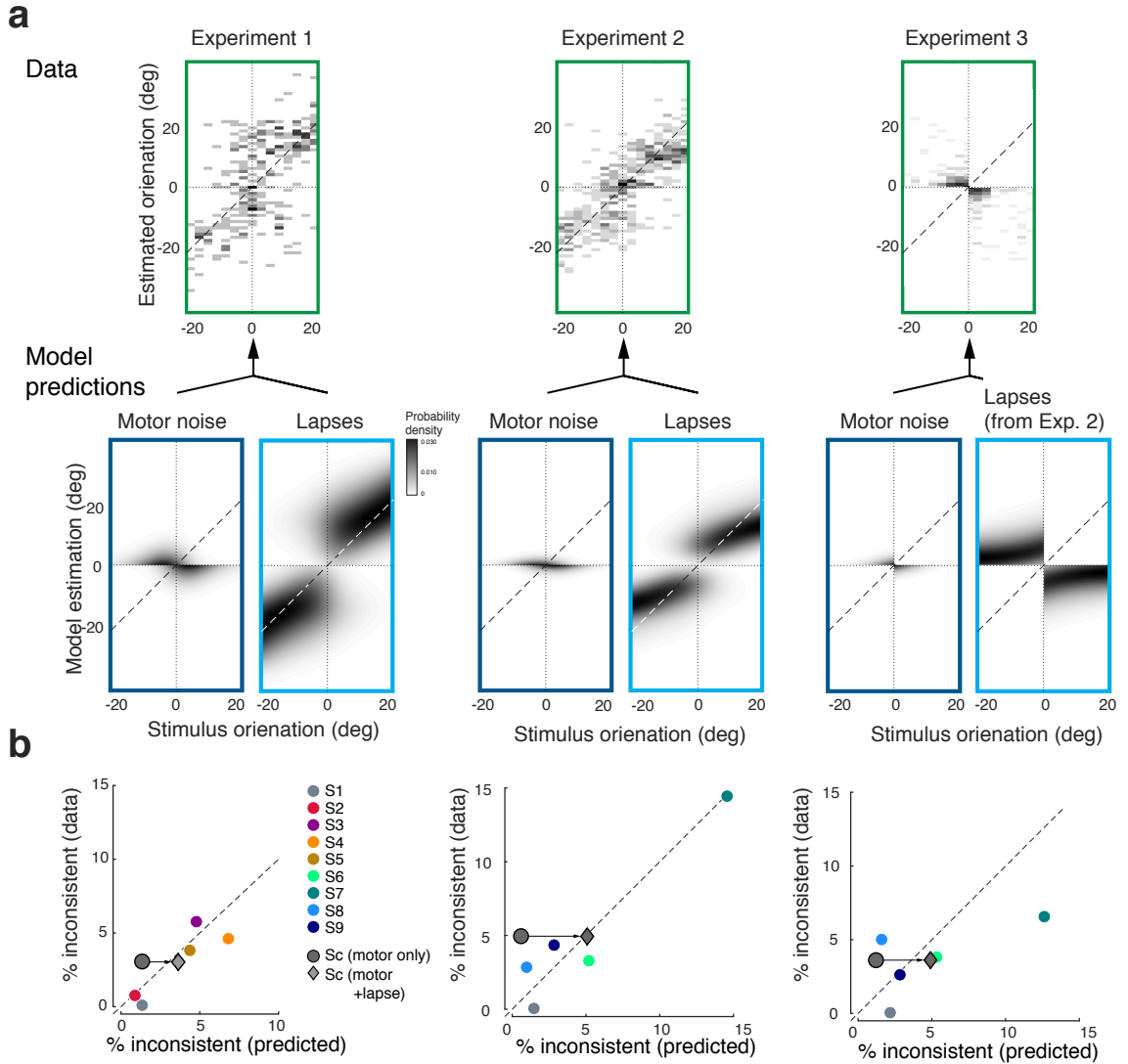


Figure 2.8: *Inconsistent trials are due to lapses and motor noise.* (a) Distribution of estimates for the small fraction of inconsistent trials (4% of the data) in each experiment (across all subjects and stimulus noise conditions). The estimation patterns can be explained as a weighted superposition of two sources of erroneous, non-perceptual behavior: lapses and motor noise. The self-consistent model well predicts the estimation patterns. All predictions are based on parameter values taken from the model fit to the consistent trial data (see Methods section 2.8). Lapse rates were extracted from the psychometric functions of the discrimination judgment for the total data. Motor noise was measured in a control experiment (see Figure supplements section 2.8, Figure supplement 2.12). (b) Quantitative predictions for each subject's total fraction of inconsistent trials are compared to the measured fractions. Predictions for the combined subjects suggest that inconsistent trials are mainly due to lapses.

judgment. For Experiments 1 and 2, the predicted estimation patterns for lapse trials are similar to the predictions for consistent trials (Figs. 2.3c, 2.5c), which makes sense since we assumed that subjects performed the estimation task correctly (*i.e.* were self-consistent) but mistakenly pushed the wrong button in the discrimination task. For Experiment 3, the pattern is different because the misremembered answer cue is always incorrect, and thus subjects estimates are based on the long tail of the sensory measurement distribution. In contrast, motor noise leads to inconsistent trials when it accidentally deflects subjects' reported orientation estimates to the other side of the discrimination boundary. Thus for all experiments they are predicted to be limited to stimulus orientations close to the boundary. A visual comparison between the measured and the predicted estimation pattern (Fig. 2.8a) confirms that the small fraction of inconsistent trials are qualitatively well explained as the combined effect of errors due to lapses and motor noise. Furthermore, we can quantitatively predict individual subjects' overall fraction of inconsistent trials based on their fit model parameter values, and measured lapse rates and motor noise (Fig. 2.8b). Analyzing the contribution of each error source further reveals that the majority of inconsistent trials are caused by lapses in the discrimination task.

2.4.3. Explaining existing experimental data

In a recent study, Zamboni et al. (2016) run different variations of the original experiment (Jazayeri and Movshon, 2007). Specifically, they manipulated the presence and orientation of the discrimination boundary at the time of the estimation task, as well as whether subjects had to explicitly perform the discrimination task or not. We fit our model to this dataset (combined subject) and the results are shown in Fig. 2.9. Experiment 1a was an exact replica of the original experiment (Jazayeri and Movshon, 2007). The observed bias patterns are consistent with the original results as well as the results from our Experiment 1, and thus well accounted for by our model (Fig. 2.9a). Experiment 1b was identical to 1a except that the discrimination boundary was removed right after subjects performed the discrimination task. This manipulation led to an increase in variance and a loss of bimodality in the distribution of estimates (Fig. 2.9b). Interestingly, however, the data

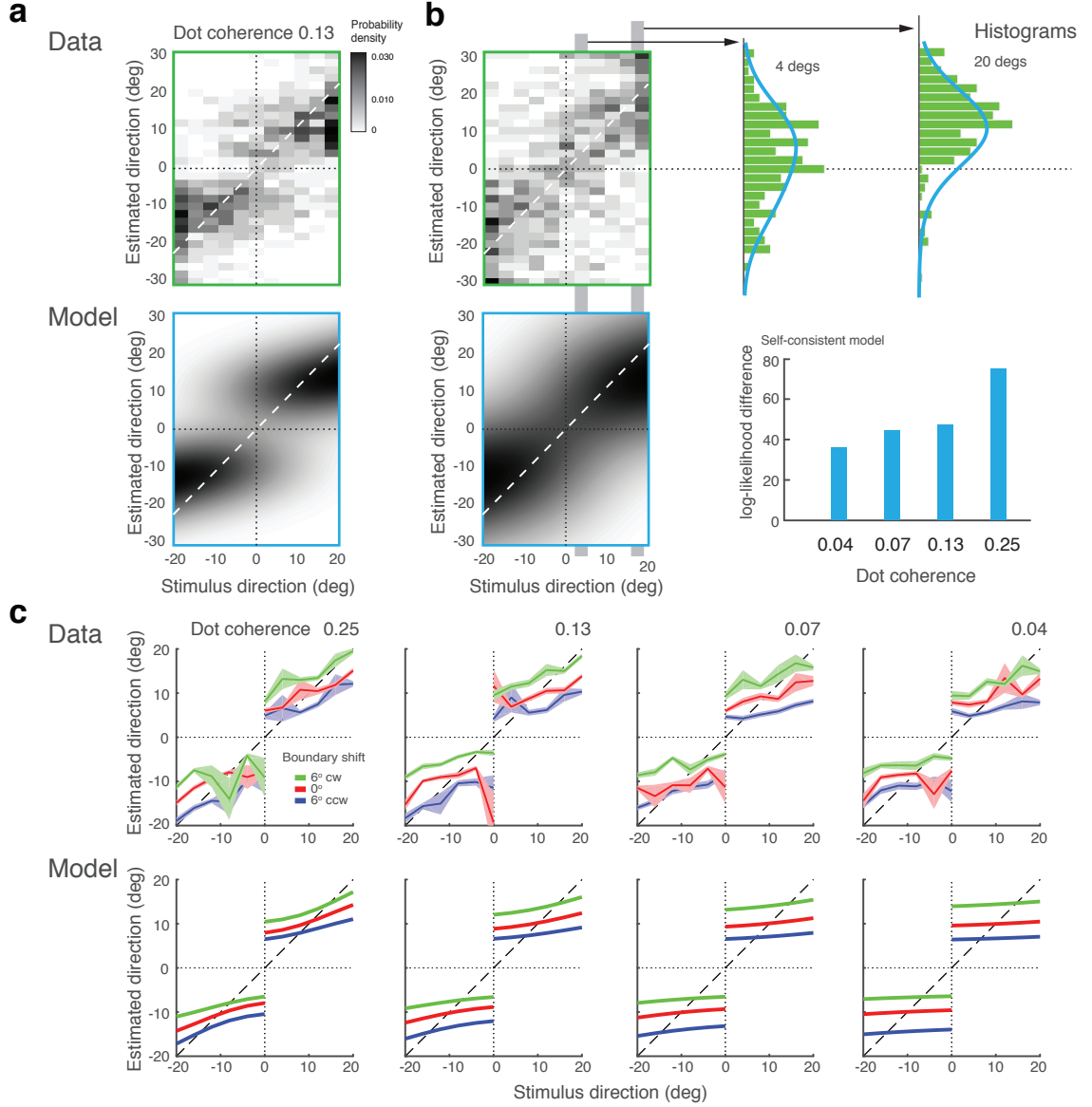


Figure 2.9: *Model fits for experimental data by Zamboni et al. (2016).* (a) Experiment 1a: Exact replication of the original experiment Jazayeri and Movshon (2007). Exemplarily shown is the estimation data (combined subject) at one stimulus coherence level (0.13) together with our model fit. (b) Experiment 1b was identical except that the boundary was not shown during the estimation task. Estimate distributions are no longer bimodal yet the self-consistent observer, relying on a noisy memory of the boundary orientation, consistently better fit the data than the strictly unimodal, independent observer model (log-likelihood difference). (c) Experiment 2 introduced a shift in the boundary orientation right before the estimation task, which subjects were not aware of (± 6 degrees). Subjects' estimates were shifted accordingly (combined subject). The self-consistent model can account for the shift if we assume that the conditional prior is applied to the shifted boundary orientation. See Fig. 2.9—figure supplement 2.18-2.20 for distributions, fits, and goodness-of-fits for all conditions.

are consistently better fit by the self-consistent model than by the independent model that strictly predicts a unimodal distribution, if we assume that the observer had to rely on a noisy memory representation of the discrimination boundary for the estimation task. A detailed inspection of the estimate distributions shows that they are wider the closer the stimulus direction is to the boundary and generally skewed towards the boundary. This suggests that subjects behaved according to the self-consistent observer model, yet the characteristic bimodal estimation pattern is hidden in the extra variance introduced by the uncertainty about the boundary orientation.

In Experiment 2, the boundary orientation was shifted by a small amount (± 6 degrees) before subjects had to perform the estimation task. Introducing a short blank screen right before the shift ensured that subjects were not aware of this manipulation. In contrast to Experiment 1, subjects were only asked to perform the estimation task. Subjects' estimates still show the same characteristic bimodal distribution although they are shifted according to the boundary shift. This suggests that subjects implicitly performed the discrimination task even though they were not asked to report an explicit judgment, which is supported by the good model fit (Fig. 2.9c). Based on these results, we conclude that self-consistent inference takes place at the time of the estimation task, can occur on memorized boundary information, and does not require an explicit discrimination judgment.

2.5. Self-consistency despite memory degradation

To what degree is self-consistent inference a necessary condition for self-consistent behavior? If working memory were perfect (*i.e.* the sensory signal m and its memory recall m_m are identical) then any reasonable observer model would be self-consistent. However, this is an unlikely scenario because it is fairly well established that continuous visual information is degrading rather quickly over time in working memory (Wilken and Ma, 2004; Bays et al., 2011). We thus expect working memory degradation to affect perceptual behavior, in particular in Experiments 1 and 2 where the average time between stimulus presentation and the estimation task was on the order of 2-3 seconds. This is supported by the model fits that revealed non-zero memory noise levels. In order to quantify how challenging working

memory noise is for maintaining self-consistency, we computed the fractions of inconsistent trials that we would expect without self-consistent inference, based on the fit memory noise levels.

Figure 2.10a shows the predicted fractions of inconsistent trials as a function of stimulus orientation for every subject and stimulus noise condition. The curves reflect the fraction of trials in which the sensory measurement m and the working memory recall m_m are on different sides of the decision boundary. Predictions vary for individual subjects yet are typically large in particular for orientations close to the boundary.

A comparison with the actually observed fractions of inconsistent trials in Experiments 1 and 2 reveals that those are much smaller and relatively independent of the stimulus orientation (combined subject for statistical reasons; Fig. 2.10b), in line with our previous conclusion that inconsistent trials predominantly reflect lapses (see Fig. 2.8). This is further supported by a comparison with subjects' individual memory noise levels: the predicted fractions are almost perfectly correlated with memory noise whereas no such correlation can be found for the observed fractions (Fig. 2.10c). Thus, above analysis suggests that if memory noise is present, the proposed self-consistency constraint is necessary in order to account for the low fractions of inconsistent trials in the data.

2.6. Alternative interpretations

While various dependencies in sequential perceptual choice tasks have been reported, such as dependencies between decision outcomes (Fernberger, 1920; Senders and Sowards, 1952; Gold et al., 2008; Fründ et al., 2014; Abrahamyan et al., 2016), decision confidence (Rahnev et al., 2015; Fleming and Daw, 2017), reaction times (Laming, 1979; Link, 1975), and error rates (Vervaeck and Boer, 1980), the impact of subjects' choices on their immediate subsequent perceptual judgments has not yet been considered a cause for sequential dependencies. Notably, Jazayeri and Movshon (2007) interpreted the reported post-decision biases as the result of a selective read-out strategy by which the brain preferentially weighs signals from those sensory neurons that are most informative with respect to the discrimination task, yet is then compelled to use the same weighted read-out signal when performing

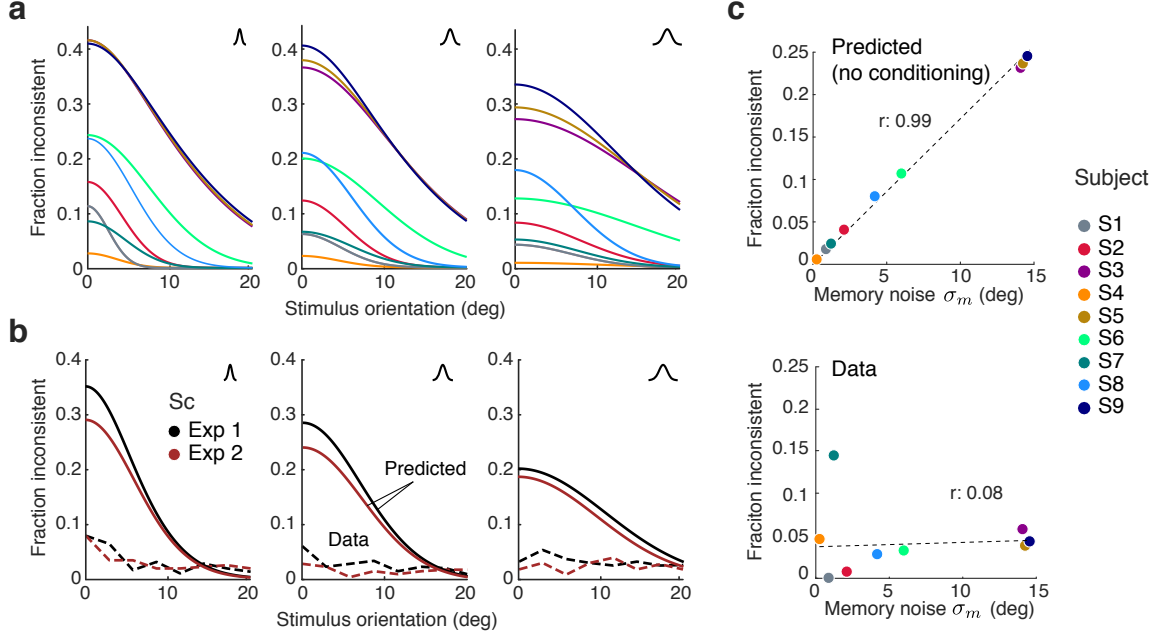


Figure 2.10: *Maintaining self-consistency in the face of working memory noise.* (a) Shown are the predicted fractions of inconsistent trials if orientation estimates are not conditioned on the preceding judgment. These are trials for which the sensory signal m and its memory recall m_m are on different sides of the discrimination boundary. Using the fit model parameters from Exp. 1 and 2, each line represents the fraction of inconsistent trials as a function of stimulus orientation for every subject (color code on the right). Each panel is for one of the three stimulus noise conditions. These large fractions are predicted for any non-trivial model whose discrimination judgment is based on m and the estimate on m_m but does not condition the estimation process on the preceding discrimination judgment. For simplicity, we did not include lapses and motor error for this analysis and thus these predictions reflect the direct consistency benefit of conditioning the estimate on the preceding discrimination judgment. (b) The actual fractions of inconsistent trials are much lower and relatively independent of stimulus orientation as they are mostly due to lapses (see Fig. 2.8b); shown is the combined subject. (c) The benefit of self-consistent inference is substantial for larger memory noise; predicted fractions are almost perfectly correlated with the fit memory noise σ_m of individual subjects. In comparison, the actual fractions of inconsistent trials are uncorrelated with memory noise levels, in line with our previous analysis showing that they are mainly due to lapses and motor noise.

the subsequent estimation task. They conclude that a non-uniform read-out profile that more strongly weighs neurons with preferred tuning slightly away from the discrimination boundary could explain the repulsive bias patterns. This is a more mechanistic, neural interpretation, which complicates a direct comparison with our normative computational model. Nevertheless, there is a fundamental difference between this interpretation and our self-consistent model in the way how it assumes the two perceptual tasks to interact: the read-out model proposes a feed-forward process where both tasks are performed independently based on the same weighted sensory signal, whereas our model assumes that feed-back of the categorical judgment is causally affecting the estimation process. Despite this difference, our experimental results alone may not be sufficient to disambiguate between these two interpretations. Experiments 2 and 3 were foremost designed to test the specific aspects of our self-consistent observer model. Although subjects' behavior in these experiments is not compatible with the originally proposed rationale for the particular shape of the read-out profile (i.e. optimized for the discrimination task - Jazayeri and Movshon (2007)), relaxing this assumption by also allowing stimulus prior information to determine the shape of the profile may lead to an alternative explanation of our data. Future work must show whether this is true or not. It seems important that such work can establish a principled description of how stimulus prior information ought to be reflected in the weighting function, otherwise the model's explanatory power will be reduced to that of a curve fit. It seems also important that any potential model assessment takes full advantage of the richness of the behavioral data, *i.e.* models should be evaluated based on their ability to account for the entire distribution of subjects' discrimination-estimation response pairs, and not only on summary statistics such as mean bias (Jazayeri and Movshon, 2007; Zamboni et al., 2016).

More challenging to reconcile with the read-out model, or any other model that does not impose some form of self-consistency constraint, are the experimental results by Zamboni et al. (2016) in combination with our consistency analysis (Fig. 2.9). The results by Zamboni et al. (2016) suggest that a different (or at least adjusted) read-out profile must be applied at

the time of the estimation task, which implies that the sensory signal up to that point needed to be stored in some form of working memory. With working memory quickly deteriorating over time (Wilken and Ma, 2004; Bays et al., 2011), our consistency analysis shows that the observed degree of trial consistency cannot be achieved by an observer model that does not condition the estimate on the discrimination judgment (Fig. 2.9). Future research is necessary to validate the levels of working memory noise we have determined with our model.

2.7. Discussion

We have shown that in a discrimination-estimation task sequence, the estimated value of a stimulus variable is systematically biased by the preceding discrimination judgment about that variable. We have introduced a self-consistent Bayesian observer model that provides an accurate and complete description of perceptual behavior in such sequential decision-making. The model assumes that the estimates are the result of a Bayesian inference process over a hierarchical generative model, which, however, is conditioned not only on the sensory evidence but also on the subject’s judgment in the preceding discrimination task. This guarantees that discrimination judgments and estimates in any given trial are consistent even when the observer must rely on working memory signals that are noisy. We show that an observer that considers the tasks independently will substantially fail to provide the level of consistency observed in the data. With a set of targeted psychophysical experiments we verified that the observed bias pattern generalizes for different low-level visual stimuli (Experiment 1), and validated the self-consistent model by showing that the pattern indeed depends on subjects’ knowledge of the stimulus prior (Experiment 2) and that subjects use their own decision as if it was correct (Experiment 3). We further validated the model with existing data from experiments that manipulated the presence and orientation of the discrimination boundary. Successful fits of the proposed observer model to individual subjects data across the various experiments demonstrate the power and accuracy of the model, and its ability to generalize across experimental conditions. Furthermore, the model fits provide a meaningful interpretation of the substantial between-subject differences

in behavior in terms of individual differences in noise levels and knowledge of the stimulus prior.

Our results and in particular the proposed self-consistent inference model have broad implications for understanding human decision-making in general. For example, subjects did not distinguish between a decision outcome they generated themselves and a decision outcome that was given to them (see section 2.4.1), which implies that for the purpose of performing the estimation task they treat their own subjective judgment as if it was correct. Computationally, this is interesting because on one hand it apparently seems to violates optimal behavior in terms of overall perceptual accuracy (obviously, a subjective judgment can be wrong). On the other hand, however, it guarantees that the observer remains self-consistent throughout the task sequence even when noise is corrupting the sensory information in working memory (Fig. 2.10). This is consistent with previous results showing that selectively discarding evidence (a seemingly irrational behavior) can improve performance when decision formation is corrupted by internal neural noise (Tsetsos et al., 2016), and thus may be rational after all. Future work is needed to investigate in more detail the impact of self-consistent inference on choice performance and perceptual accuracy.

2.8. Method and supplementary material

Methods

Experimental setup

Ten subjects with normal or correct-to-normal vision (6 males, 4 females; one non-naïve) participated in the experiments. One subject (male) was excluded from the analysis because he failed to correctly execute the estimation task. All subjects provided informed consent. The experiments were approved by the Institutional Review Board of the University of Pennsylvania under protocol #819634.

General methods: Subjects sat in a dimmed room in front of a special purpose computer monitor (VIEWPixx3D, refresh rate of 120 Hz and resolution of 1920 x 1080 pixels). Viewing distance was 83.5 cm and enforced with a chin rest. We programmed all experiments

in Matlab (Mathworks, Inc.) using the MGL toolbox (<http://justingardner.net/mgl>) for stimulus generation and presentation. The code was run on an Apple Mac Pro computer with Quad-Core Intel Xeon 2.93 GHz, 8GB RAM. Subjects were asked to fixate at the fixation dot whenever it appeared on the screen. Before subjects did the main experiments, they had 2-3 training sessions during which they familiarized themselves with the discrimination and the estimation task. After that, every subject either completed 1800 trials in 3-4 sessions for Experiment 1 or completed 3600 trials in 6-8 sessions for Experiment 2 and 3. This is equivalent to 40 trials per every one of the 15 stimulus orientations and the three noise conditions. Sessions lasted approximately 50 minutes. Subjects used a gamepad (Sony PS4 Dualshock) as input device: they reported their decision in the discrimination task by pressing the appropriate trigger button (left for 'ccw', right for 'cw'), and used the analog joystick of the gamepad to indicate their perceived stimulus orientation by adjusting a reference line and then confirming their estimate with an additional button press. Screen background luminance was 40 cd/m² and mean stimulus luminance was 49 cd/m².

Experiment 1: Five subjects (S1-5) participated in Experiment 1. In each trial, subjects viewed a white fixation dot (diameter: 0.3°) and two black marks (length: 3°, distance from fixation: 3.5°) indicating a discrimination boundary whose orientation was randomly chosen around the circle. After 1300 ms, the orientation stimulus consisting of an array of white line segments (length: 0.6°) was presented for 500 ms. The array consisted of two concentric circles of lines centered on the fixation: the outer (diameter: 3.8°) contained 16 line segments and the inner (diameter: 1.8°) contained 8 line segments. Small random variations (from -0.15° to 0.15°) were independently added to the x-y coordinates of each line segment. The orientation of every line segment was drawn from a Gaussian distribution with mean given as one of 15 stimulus orientations relative to the boundary (from -21° to 21° in steps of 3°) and standard deviation σ as one of 3 values (0°, 6° and 18°). After the stimulus disappeared, subjects were asked to indicate whether the overall orientation of the array was clockwise or counter-clockwise relative to the boundary. If subjects responded within 4 seconds, they then were also asked to indicate their perceived stimulus orientation.

Otherwise the current trial was skipped and added to the back of the trial list. Every trial was followed by a randomly chosen inter-trial interval of 300 ms to 600 ms duration (blank screen; mean luminance).

Experiment 2: Five subjects (S1 and S6-9) participated in Experiment 2. The procedure was identical to Experiment 1 except that at the beginning of every trial, subjects were reminded of the stimulus distribution by presenting a prior cue consisting of a gray arc for 800 ms. The arc (width: 0.2° , eccentricity from fixation: 3.5°) spanned the range $\pm 21^\circ$ relative to the discrimination boundary indicating the total true stimulus distribution. Subjects were instructed that stimulus orientation was guaranteed to occur anywhere within this range with equal probability.

Experiment 3: The same five subjects that participated in Experiment 2 also participated in Experiment 3. The procedure was identical to Experiment 2 except for the following: First, the prior cue was present only for 500 ms, after which it was reduced to a colored arc that only spanned the orientation range at the side of the discrimination boundary where the stimulus orientation in this trial would occur. This colored cue indicated the correct answer ('cw' or 'ccw') and was shown for 300ms. Its color (red or green) was randomly assigned and uncorrelated with the stimulus orientation. Second, instead of the orientation discrimination task, subjects were tasked to recall the color of the cue.

Control Experiment (Motor noise): As part of the training, all subjects participated in a control experiment that allowed us to estimate subjects' individual motor noise levels. They were first presented with a fixation dot and boundary black marks (like in Experiment 1-3). After that, they were presented with an orientation stimulus consisting of a single white line (like the reference line in Experiment 1-3, see *e.g.* Fig. 2.1a) for 500 ms. They then had to reproduce the stimulus by adjusting said reference line with the analog joystick of the gamepad. After subjects pressed a button to confirm their response, they received feedback in form of a green reference line indicating the true stimulus orientation. This control experiment consisted of 600 trials. The boundary orientation was uniformly sampled around the circle. The stimulus orientation was uniformly distributed around the boundary

as in Experiment 1-3. We used the measured standard deviation σ_0 in subjects' estimates as a direct measure of subjects' individual motor noise levels. Fig. 2.3—figure supplement 2.12 shows the measured σ_0 for all 9 subjects. We used these measured levels as fixed parameters in all our model fits and predictions, assuming that motor noise is Gaussian and is uniform and independent of the reproduced orientation relative to the discrimination boundary.

Self-consistent Bayesian observer model

The model is formulated as an observer that performs Bayesian inference over the hierarchical generative model shown in Fig. 2.2a. The observer is assumed to solve the two perceptual inference tasks in sequence: the observer first judges whether the stimulus orientation was clockwise (cw) or counter-clockwise (ccw) of a random discrimination boundary, and then performs an estimate of the actual stimulus orientation. The key feature of the model is that the inference process associated with the estimate is *conditioned* on the observer's preceding categorical judgment. As such, the observer treats their own, subjective categorical judgment as if it was a correct statement about the stimulus (see Experiment 3). In the following we describe the Bayesian formalism of this conditioned task sequence.

Discrimination judgment: Let θ be the true stimulus orientation relative to the discrimination boundary, m the noisy sensory measurement of the stimulus orientation at presentation, and $C = \{\text{'cw'}, \text{'ccw'}\}$ the categorical decision variable indicating whether the stimulus orientation is clockwise (cw) or counter-clockwise (ccw) of the discrimination boundary. Assuming a symmetric loss structure (*i.e.* no benefit for one correct decision over the other), the observer solves the categorical decision task by picking the category with maximal posterior probability given the sensory measurement m , thus

$$\hat{C}(m) = \operatorname{argmax}_{C \in \{\text{'cw'}, \text{'ccw'}\}} p(C|m) . \quad (2.1)$$

The decision process $\hat{C}(m)$ is a deterministic mapping from any particular m to either one

of the two categories. The posterior distribution $p(C|m)$ is given as

$$p(C|m) \propto p(m|C)p(C) , \quad (2.2)$$

where $p(C)$ is the prior probability and $p(m|C)$ the likelihood over the choices. We can obtain this likelihood by marginalizing the stimulus likelihood over all stimulus orientations, that is

$$p(m|C) = \int_{-\pi}^{\pi} p(m|\theta)p(\theta|C)d\theta . \quad (2.3)$$

The stimulus likelihood is fully determined by the noise characteristics of the sensory measurement m , thus by the distribution $p(m|\theta)$ of sensory measurements m for any given true stimulus orientation θ . Finally, the model's prediction of the psychometric function in the decision task is obtained by marginalizing the posterior (Eq. (2.2)) for a particular choice (*e.g.* "Fraction reported cw" - see Fig. 2.2b) over the sensory measurement distribution,

$$p(\hat{C} = \text{'cw'}|\theta) = \int p(\hat{C} = \text{'cw'}|m)p(m|\theta)dm . \quad (2.4)$$

Estimation task: Following the categorical judgment, the observer then solves the estimation task by computing the mean of the posterior distribution (*i.e.* assuming a loss function that minimizes squared error) over θ . In contrast to the independent observer, however, we assume that the posterior probability is *conditioned on the observer's own choice* \hat{C} in the preceding categorical decision task. Because the stimulus has long disappeared by the time the subject performs the estimation task (see experimental design), we formulate the posterior on a memorized version of the sensory measurement. With m_m referring to a noisy recall sample from working memory of the measurement m (doubly stochastic) we write the optimal conditioned estimate as

$$\hat{\theta}(m_m, \hat{C}) = \int_{-\pi}^{\pi} \theta p(\theta|m_m, \hat{C})d\theta , \quad (2.5)$$

with the posterior

$$p(\theta|m_m, \hat{C}) = \frac{p(m_m|\theta, \hat{C})p(\theta|\hat{C})}{p(m_m|\hat{C})} . \quad (2.6)$$

The estimate $\hat{\theta}(m_m, \hat{C})$, even though the result of Bayesian inference, describes a deterministic mapping from any particular m_m to an estimate. However, we obtain two different formulations for the estimate, one for each potential categorical judgment³.

With Eq. (2.5) the distribution of the estimates can be computed based on the distribution of the memorized sensory measurement

$$p(m_m|\theta, \hat{C}) = \int p(m_m|m)p(m|\theta, \hat{C}(m))dm . \quad (2.7)$$

Note that the above marginalization is limited to measurements m that led to the particular categorical judgment $\hat{C}(m)$ (as given by Eq. (2.2)).

The model's description of the conditioned distribution of the estimates $p(\hat{\theta}|\theta, \hat{C})$ is obtained by a variable transformation for the conditional measurement distribution $p(m_m|\theta, \hat{C})$, *i.e.* substituting m_m with the estimate $\hat{\theta}(m_m, \hat{C})$ (Eq. (2.5)). Finally, the model's prediction of the entire distribution of the estimates $p(\hat{\theta}|\theta)$ (*i.e.* the density plots shown in *e.g.* Fig. 2.2c) is given by marginalizing over the decision outcomes, thus

$$p(\hat{\theta}|\theta) = \sum_{\hat{C}} p(\hat{\theta}|\theta, \hat{C})p(\hat{C}|\theta) . \quad (2.8)$$

Estimation task with known correct category assignment (Experiment 3):

If the observer knows the category assignment upfront (as in Experiment 3), the above formulation of solving the estimation task slightly changes in that Eqs. (2.5)–(2.8) are conditioned on the actual correct category assignment C rather than a the inferred category \hat{C} . In particular, this changes the marginalization over m in Eq. (2.7) such that it is no longer limited to values of m that are consistent with a desired category assignment (according to $\hat{C}(m)$), and also Eq. (2.8) where the sum is over the actual binary category

³In general, there are as many different versions for the estimate as there are potential choices in the categorical judgment task.

probability $p(C|\theta)$ rather than the inferred probability $p(\hat{C}|\theta)$. As a result, the predicted biases with identical model parameters are slightly smaller when the observer knows the correct category assignment compared to when the category has to be inferred first.

Specific assumptions defining the generative model: We made the following specific assumptions in defining the components of the generative model (Fig. 2.2a):

- We set the category prior $p(C) = 0.5$ for both values of C since the two choices are equally likely in all our experiments.
- The categorical stimulus prior $p(\theta|C)$ was assumed to reflect subjects' individual expectations about the experimental distribution of stimulus orientations. We modeled this prior to be identical but mirrored around the discrimination boundary for the two choices. More specifically, we assumed it to be uniform over the angle α relative to the boundary with a smooth cosine fall-off from the uniform density value to zero over the additional angle β . We then defined the prior width w_p (see Figs. 2.4b, 2.7c) as the total angle relative to the boundary where the prior density decreased to half of its uniform value, that is $w_p = (\alpha + 2/3\beta)$.
- We assumed the sensory measurements m to reflect noisy samples of the true stimulus orientation θ , with $p(m|\theta)$ to be a Gaussian with mean θ and standard deviation σ_s that is monotonically depending on the array distribution width σ of the orientation stimulus. Although σ_s was assumed to be subject dependent and thus a free parameter, we assumed that across experiments σ_s was the same for a given subject and stimulus noise condition.
- We treat the sensory evidence in the estimation task m_m as being a sample from a Gaussian with mean m (original sensory measurement) and standard deviation σ_m . We assume that σ_m is independent of stimulus uncertainty yet can be individually different for different subjects.

Independent Bayesian observer model

The “independent” observer model, as defined in our paper, is formulated on the same generative model as the self-consistent observer model (Fig. 2.3a) and thus has identical model parameters. The only formal differences are that

- the posterior over stimulus orientation is not conditioned on the discrimination judgment \hat{C} (as in Eq. (2.6)), thus

$$p(\theta|m_m) = \frac{p(m_m|\theta)p(\theta)}{p(m_m)} \quad \text{with} \quad p(\theta) = \sum_C p(\theta|C)p(C) , \quad (2.9)$$

- marginalization over the memorized sensory signal is not limited to measurements that are in agreement with a particular discrimination judgment \hat{C} (as in Eq. (2.7)), thus

$$p(m_m|\theta) = \int p(m_m|m)p(m|\theta)dm . \quad (2.10)$$

Having the same parameters as the self-consistent observer model allows a direct log-likelihood comparison in judging the goodness-of-fit.

Model fits

We jointly fit the model to the data of both the decision and estimation task by maximizing the likelihood of the model given the data:

$$p(D|\rho) = \prod_{i=1}^n P(D_i|\rho) = \prod_{i=1}^n P(\hat{C}_i|\rho, \theta)p(\hat{\theta}_i|\hat{C}_i, \rho, \theta) \quad (2.11)$$

where D is the data, ρ represents the parameters of the model, θ is the true orientation, \hat{C}_i is the decision outcome, $\hat{\theta}_i$ is the orientation estimate, i is the trial index and n is the number of trials.

For all fits, we only excluded trials with inconsistent response pairs (*i.e.* trial in which subjects discrimination judgment and estimate were not consistent in terms of their categorical identity; approximately 4% of the trials.). As we demonstrate (Fig. 2.8), the fractions

and the bias characteristics of these inconsistent trials can be fully predicted based on the fit model parameters to the consistent data by assuming that they are caused by motor noise and lapses (Fig. 2.8).

Subjects' motor noise was accounted for by assuming that the recorded orientation estimates follow the distributions of estimates $p(\hat{\theta}|\theta)$ (Eq. (2.8)) convolved with the motor noise kernel. We assume motor noise to be Gaussian with a standard deviation σ_0 that was individually determined for each subject from the control experiment (see above). Motor noise levels across subjects were fairly similar with an average $\sigma_0 = 4.3$ degrees. Fig. 2.3–figure supplement 2.12 shows measured noise levels for all subjects.

Our model fit contained a total of 6 free parameters:

- standard deviations σ_s for the 3 noise levels of the stimuli (additive Gaussian noise).
- standard deviation σ_m for memory noise (additive Gaussian).
- two parameters α, β for the prior distribution over stimulus orientation, defining its uniform range and smoothness, respectively.

The Nelder-Mead simplex algorithm was used to minimize the term $-\log(p(D|\rho))$. Twenty iterations of the optimization procedure were performed using randomized initial parameter values in order to obtain the best fitting model.

Model fit to data by Zamboni et al. (2016)

Experiment 1: For the condition where the decision boundary was always present, we fit the data with exactly the same model assumptions as we used for fitting the data from our Experiment 1. We added motor noise as a free parameter because Zamboni et al. (2016) did not use a control experiment to determine the motor noise. For the condition where the decision boundary was removed after subjects did the discrimination task, we assumed that the observer had to rely on a noisy memory recall of the true boundary orientation θ_b when performing the estimation task. As a result, the conditioned prior varies in every trial depending on that memory recall. We assumed the recalled orientation to be Gaussian distributed around the true boundary orientation with a standard deviation σ_b that was a

free parameter. Because the same group of subjects run both versions of the experiments, we first fit the self-consistent model to the data from the boundary-present condition and then used those parameters to fit the data in the boundary-absent condition with the only free parameter being σ_b (Fig. 2.9–figure supplement 2.18). When computing the goodness-of-fit (Fig. 2.9–figure supplement 2.20) we assumed the independent observer model to have the same additional noise parameter σ_b .

Experiment 2: Although the discrimination boundary was present throughout the entire trial, subjects were only asked to perform the estimation task. Furthermore, unknown to the subjects, the decision boundary was either kept the same or shifted 6 degrees (cw or ccw). For our model fit, we assumed that subjects, even not explicitly asked to do so, implicitly performed the discrimination task and then subsequently conditioned the estimation process on that implicit decision as described by the proposed self-consistent observer model. Thus we fit the data with exactly the same model assumptions used to fit the data from our Experiment 1, with the addition that for the conditions where the decision boundary was shifted, the conditioned prior was shifted accordingly. The fits shown in Fig. 2.9, Fig. 2.9–figure supplement 2.19 are joint fits to data from all three shifted conditions.

Code availability

Computer code (MATLAB) providing model simulations for all three experiments is freely available at <https://github.com/cpc-lab-stocker/Self-consistent-model> (Luu and Stocker, 2018).

Figure supplements

Experiment 1

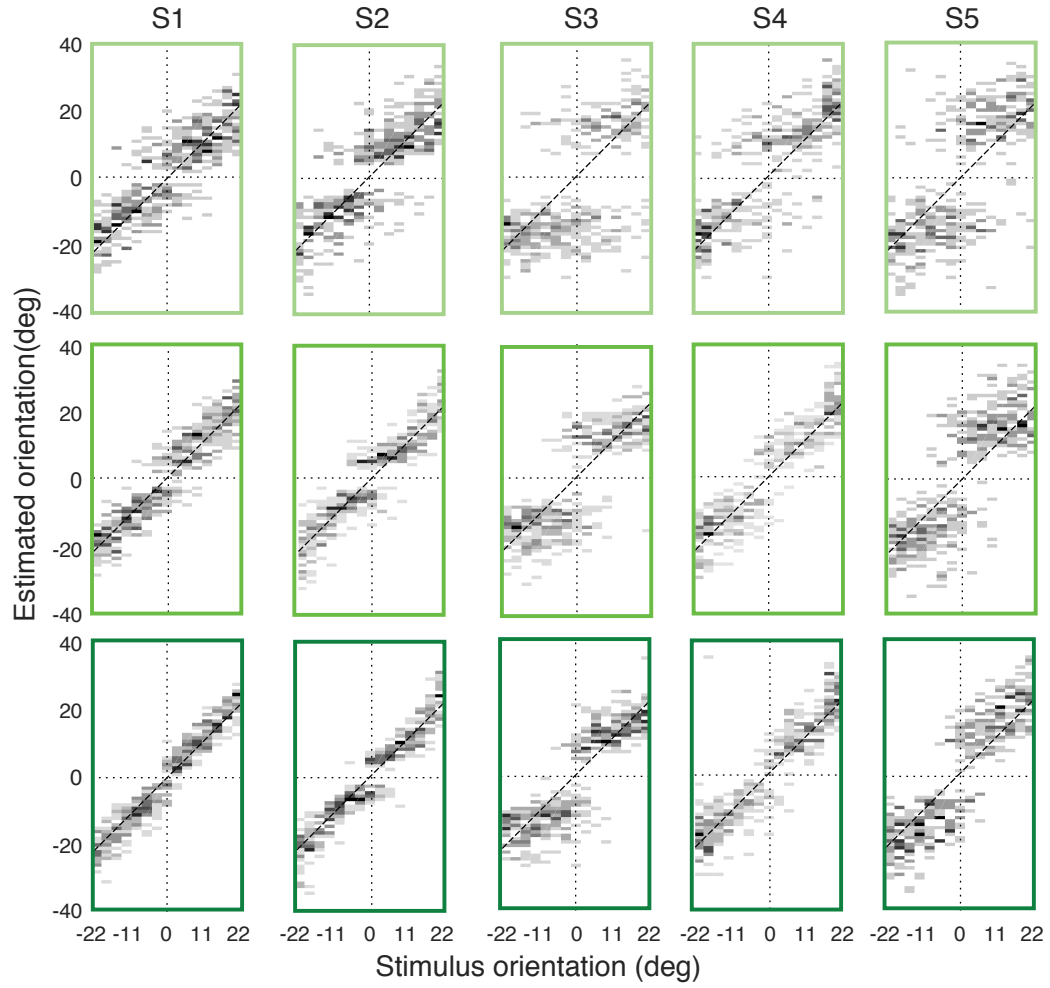


Figure 2.11: *Full distributions of individual subjects' estimates in Experiment 1.* Each row corresponds to one of the three stimulus noise conditions (color-code as in main text).

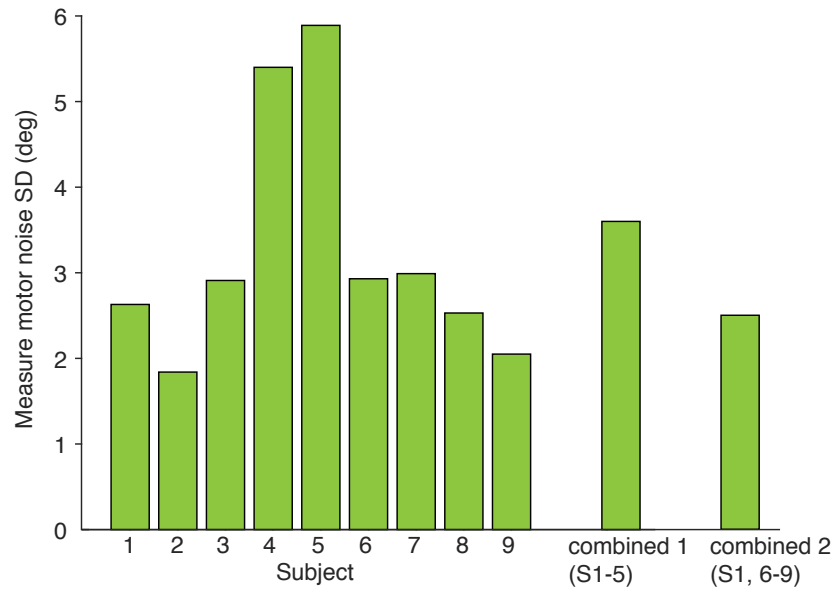


Figure 2.12: *Measured motor noise of individual subjects.* Shown are the extracted values of the standard deviation σ_0 in subjects' estimates in the control experiment (see Methods). These individual values were used in modeling the data from Experiment 1-3, assuming that motor noise was Gaussian distributed.

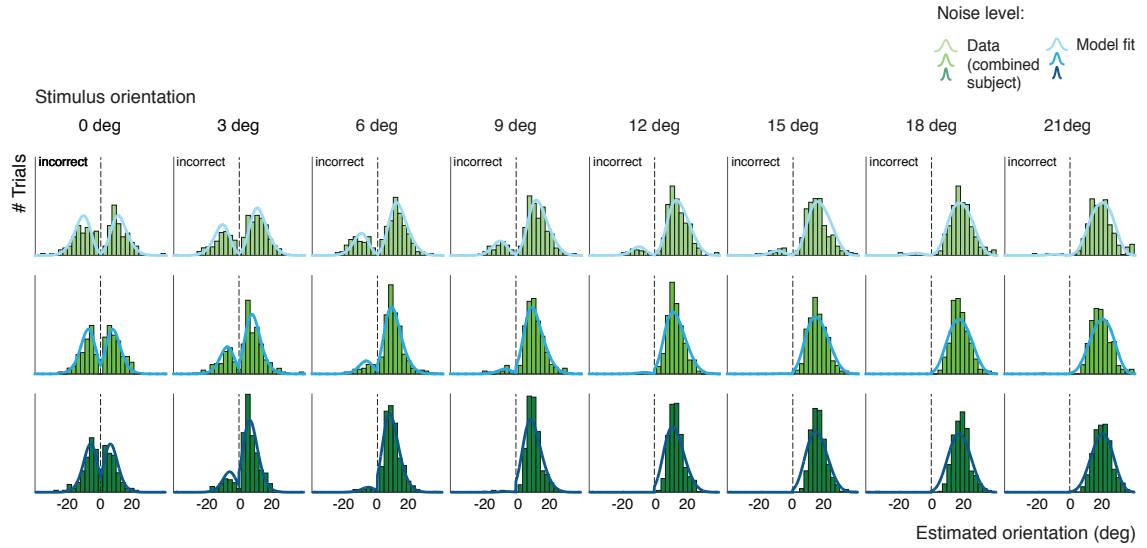


Figure 2.13: *Histogram plots of the orientation estimates together with the model fit for Experiment 1 (combined subject).* Each row is for one of the three stimulus noise conditions. Note, this figure corresponds to Fig. 2.2c in the main text.

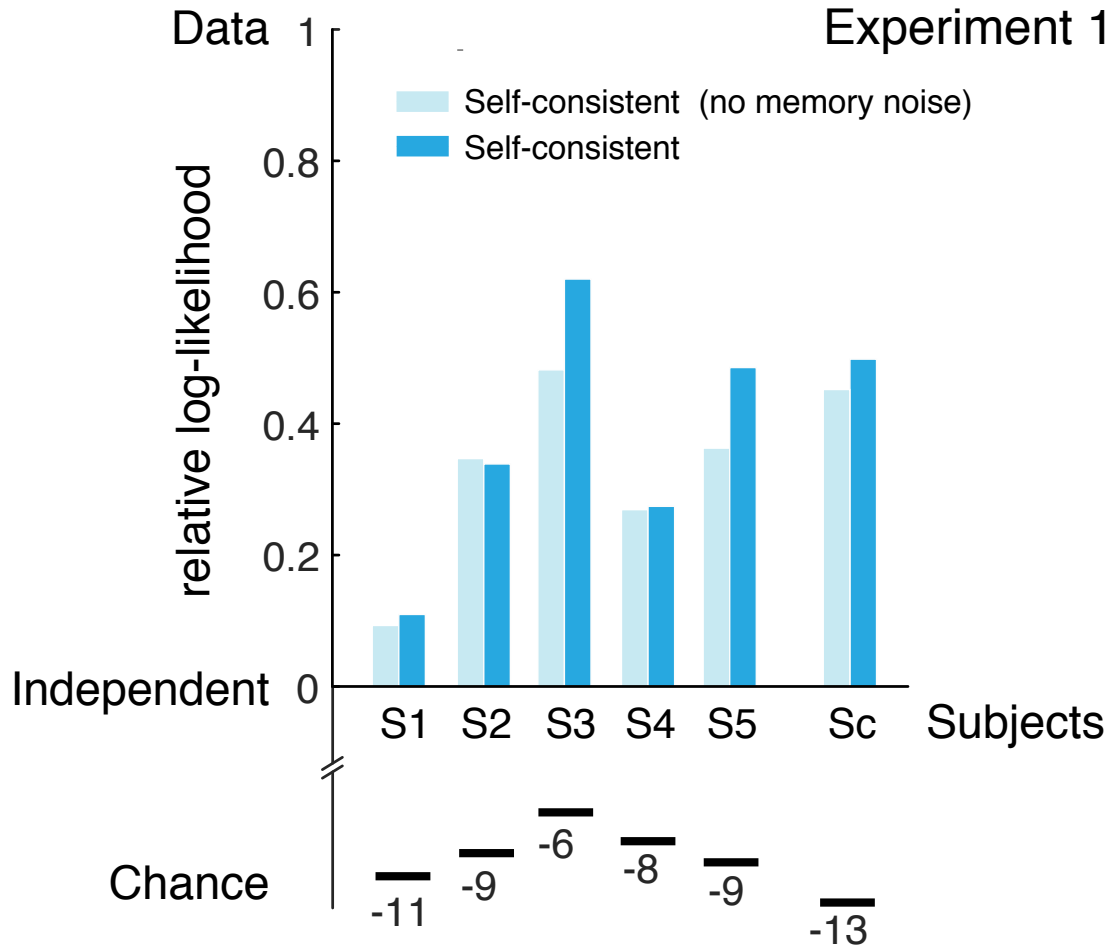


Figure 2.14: *Goodness-of-fits for Experiment 1.* Log-likelihood values of the fit self-consistent observer model for every subject (as well as the combined subject Sc), relative to the range defined by the likelihoods of the independent Bayesian observer and a hypothetical, omniscient model ('Data'). The latter can be thought of as the data explaining itself, *i.e.* a model 'defined' by the empirical probabilities of the data. The log-likelihoods of a random observer ('Chance') are also given as additional reference. This observer can be thought of as 'being blind', thus providing random answers in both the discrimination task and the estimation task (sampling from a uniform distribution). The self-consistent observer model is consistently outperforming the independent Bayesian model in explaining the data. Note, the self-consistent and the independent observer model have exactly the same model parameters. Also, a version of the self-consistent observer model that does not include noise in the memory recall of the sensory signal (Stocker and Simoncelli, 2007) generally does not fit the data as well.

Experiment 2

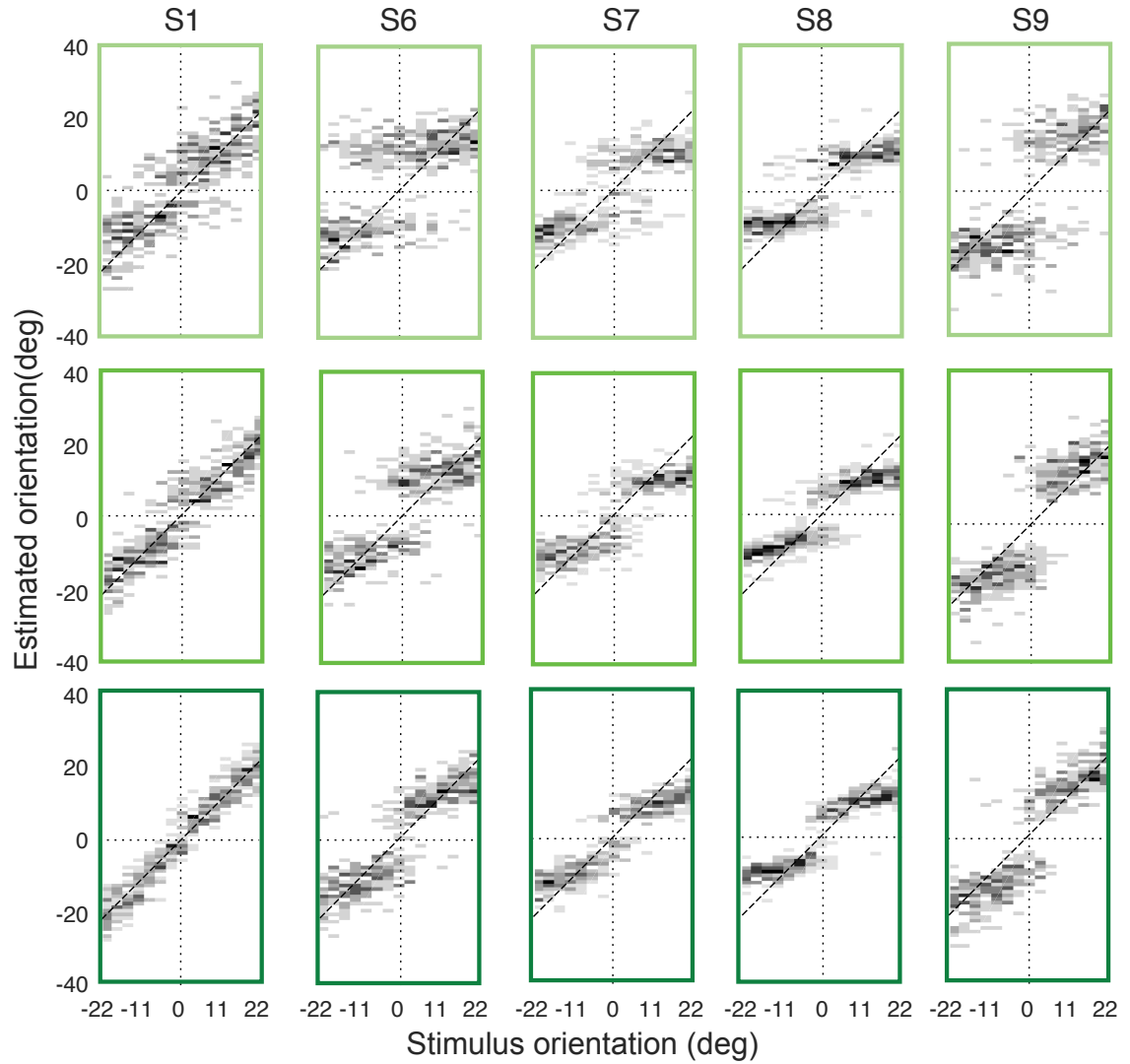


Figure 2.15: *Full distributions of individual subjects' estimates in Experiment 2.* Each row corresponds to one of the three stimulus noise conditions (top-bottom: highest-lowest; color-code as in main text).

Experiment 3

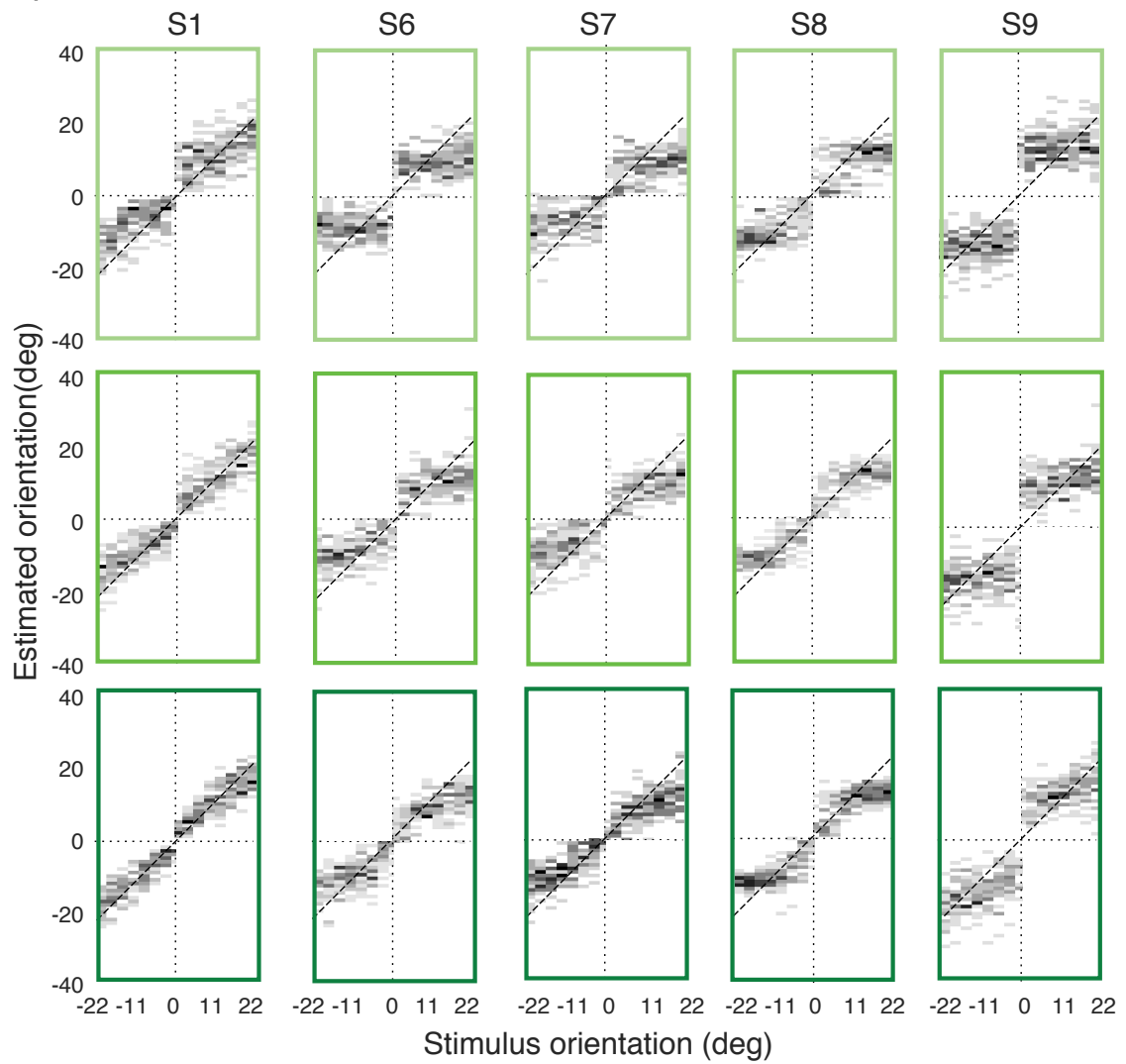


Figure 2.16: *Full distributions of individual subjects' estimates in Experiment 3. Each row corresponds to one of the three stimulus noise conditions (top-bottom: highest-lowest; color-code as in main text).*

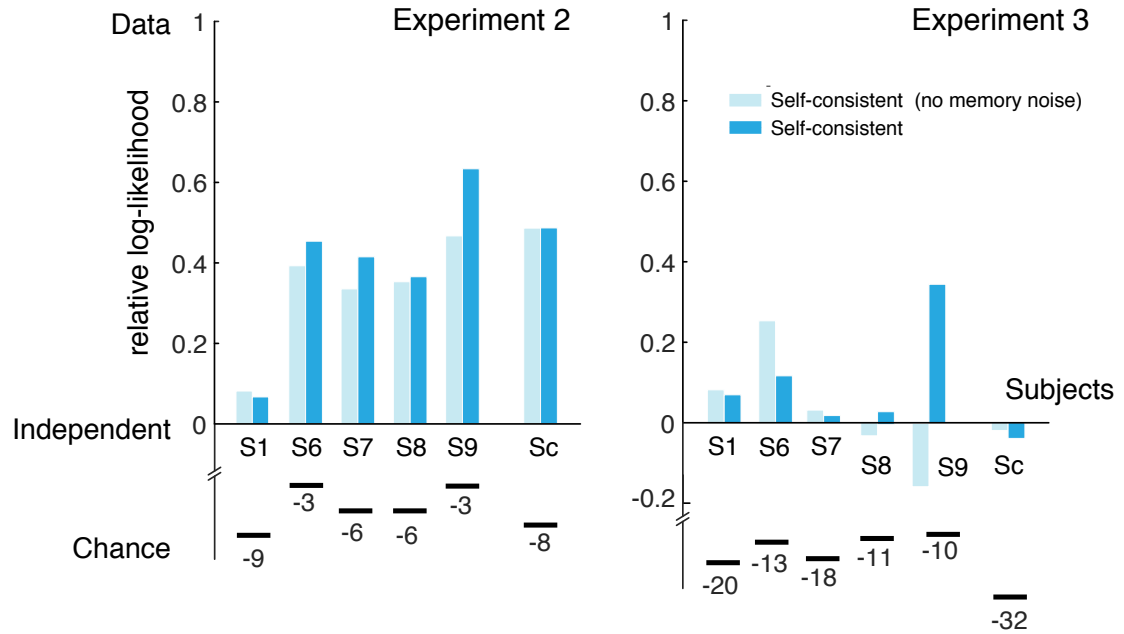
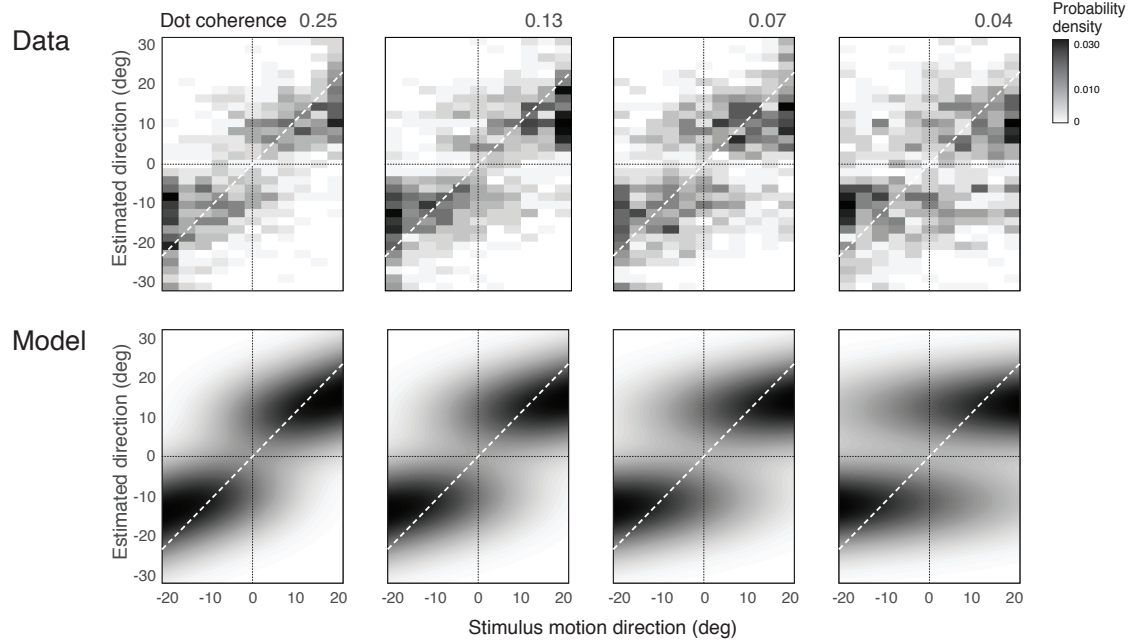


Figure 2.17: *Goodness-of-fits for Experiments 2 and 3.* Relative log-likelihood values of the fit self-consistent observer model for every subject (as well as the combined subject Sc). Relative scale is defined as described for Fig. 2.4–figure supplement 2.14. The self-consistent observer model is consistently outperforming the independent Bayesian model in explaining Experiment 2. For Experiment 3 both models are formally identical; the marginal differences in likelihood are simply because their fit parameter values slightly differ because they are from joint fits to data from both Experiment 2 and 3 (subject S1; joint fit to all three experiments). Note, the self-consistent and the independent observer model have exactly the same model parameters. Also, a model that does not include noise in the memory recall of the sensory signal generally does not fit the data as well as the full self-consistent observer model.

Boundary present:



Boundary absent:

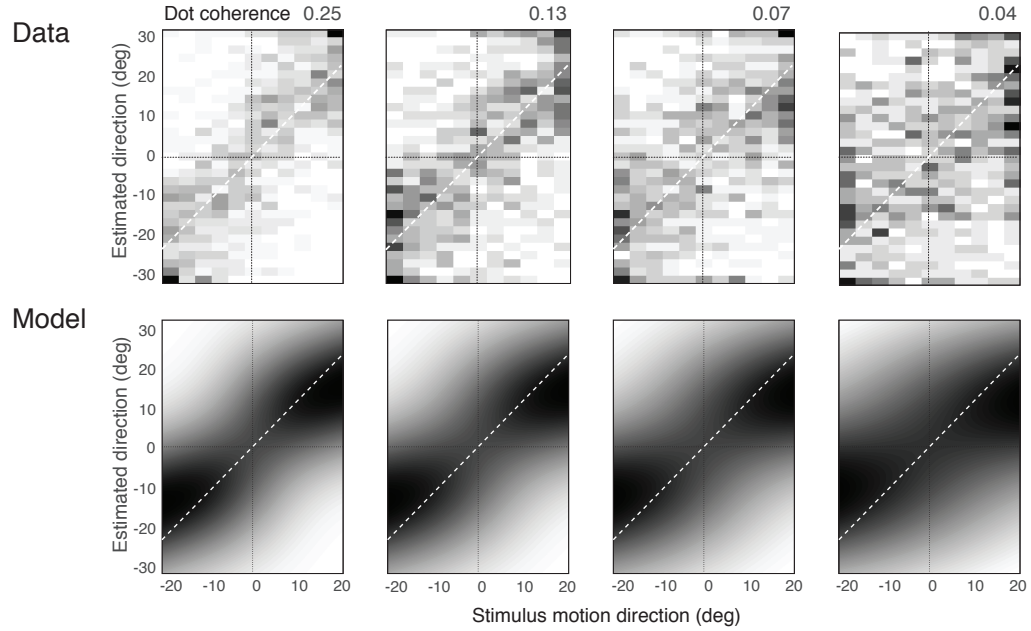


Figure 2.18: *Zamboni et al. (2016)* data (*Experiment 1, combined subject*) and fit with the *self-consistent observer model*.

Boundary shift:

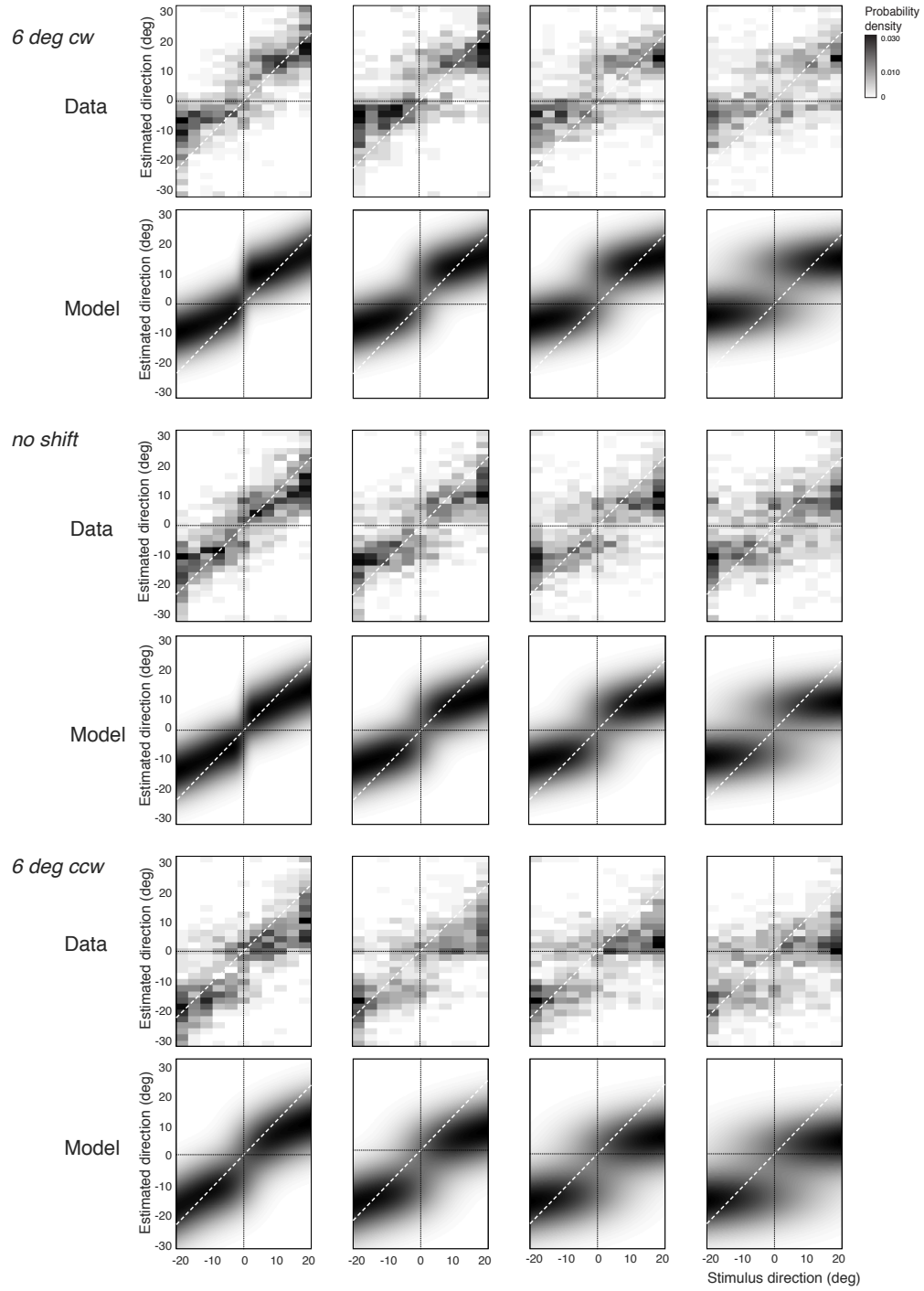


Figure 2.19: *Zamboni et al. (2016) data (Experiment 2, combined subject) and fit with the self-consistent observer model.*

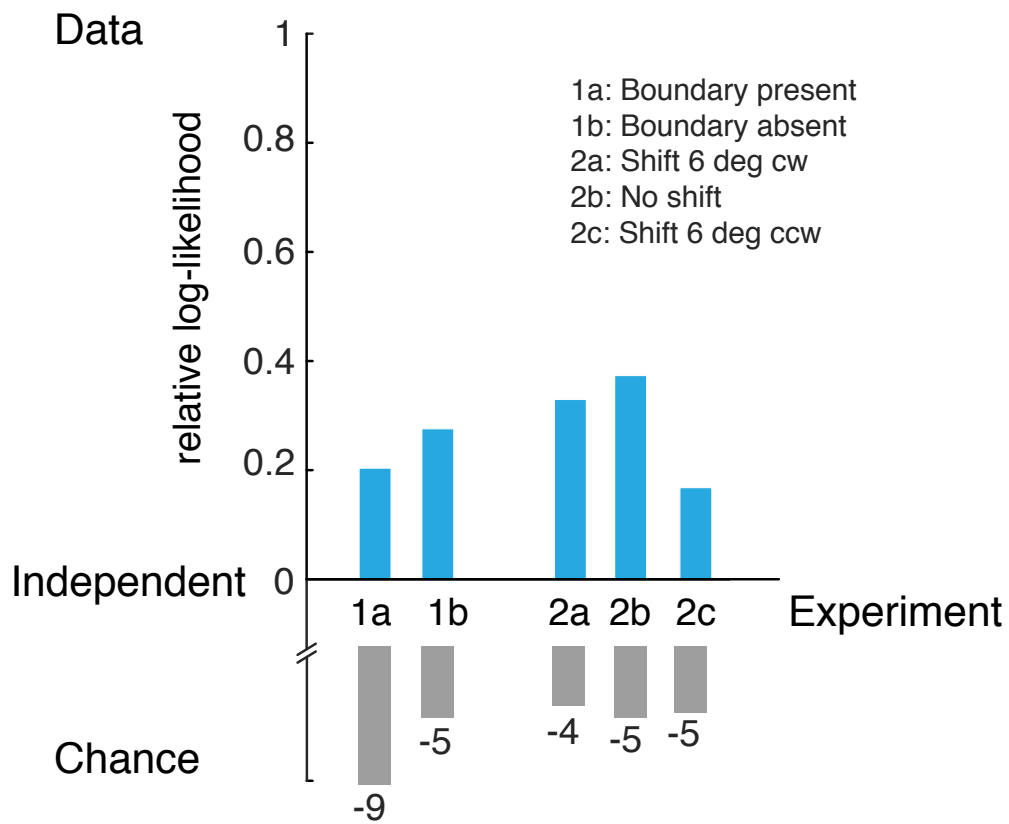


Figure 2.20: *Relative log-likelihoods of model fits for Zamboni et al. (2016) data.* Relative log-likelihood values of the self-consistent observer model fit to the combined subject data for each experiment. Relative scale is defined as described for Fig. 2.4–figure supplement 2.14.

Chapter 3

Post-decision sensory representation

3.1. Introduction

I have shown in the previous chapter that human percept of stimulus orientation is systematically biased after having made a categorical decision. The biases were well explained by a self-consistent Bayesian observer model that downweighs the sensory information inconsistent with the preceding decision. However, it is not clear what is the mechanism underlying the self-consistent inference. One hypothesis is that the sensory representation is directly modified such that the inconsistent sensory information is completely erased. That is congruent with several studies showing that making a decision may affect the subsequent accumulation of sensory evidence (Bronfman et al., 2015; Yu et al., 2015b). Alternatively, subjects may maintain the full sensory representation and the self-consistent conditioning happens at the interpretation/decoding of the sensory information. That is congruent with a study by Peters et al. (2017) that suggests although subjects show self-consistent behavior, both consistent and inconsistent sensory information are still available in the brain. Along the same line, several other studies on confidence judgment suggest that sensory evidence is preserved after a decision (Pleskac and Busemeyer, 2010; Meyniel et al., 2015).

In this chapter, I conducted two psychophysics experiments and compared alternative observer models to test the two hypotheses about the sensory representation in working memory (Fig. 3.1a). To probe the post-decision sensory representation, I provided feedback to subjects on the decision task before they performed subsequent estimation task. The results suggest that subjects maintain the full sensory representation of the stimulus. Interestingly, model analysis indicates that subjects seemed to revise the original sensory information to make it consistent with the received feedback. That demonstrates the strong

tendency of the brain to maintain self-consistency along the inferential process.

3.2. Probing post-decision sensory representation

3.2.1. Experimental test and alternative models

In Experiment 1, after being presented with an orientation stimulus, subjects had to indicate whether the stimulus orientation was clockwise (cw) or counter-clockwise (ccw) relative to a decision boundary (Fig. 3.1b). After 300 ms, a feedback tone (100% valid) was played to specify whether subjects' decision was correct or incorrect. Then subjects had to report the perceived stimulus orientation, taking into account the feedback. The experiment was similar to Experiment 1 in chapter 2 except that in the current study, subjects were given feedback in the decision task and they received a substantial bonus for good performance in both tasks.

To distinguish between different hypotheses about post-decision sensory representation, I focused on the critical trials when subjects' decision was incorrect. In hypothesis 1 (Fig. 3.2a), after making a categorical decision (e.g. 'cw'), the observer erases the part of sensory representation that is inconsistent with that decision (the ccw portion). When receiving feedback specifying that the decision is incorrect, the observer updates the prior on the stimulus range to make it consistent with the feedback. Here I consider two strategies the observer may employ to perform the estimation with the modified sensory representation. For the first strategy (1a-Prior Only), because the observer loses the sensory information on the correct side, it disregards the remaining sensory information and uses the updated prior to make an estimation. Specifically, the estimate is the mean of the updated conditional prior. For the second strategy (1b-Flip Estimate), the observer pretends that the decision is correct and estimates with the sensory information on the incorrect side. The estimate is then flipped to the correct side.

In hypothesis 2 (Fig. 3.2b), the observer maintains the full sensory representation after the categorical decision. Here I also consider two possible strategies the observer may use to perform the estimation. Both strategies assume the observer updates the conditional prior based on the feedback. The first strategy (2a-Flip Decision) assumes that the ob-

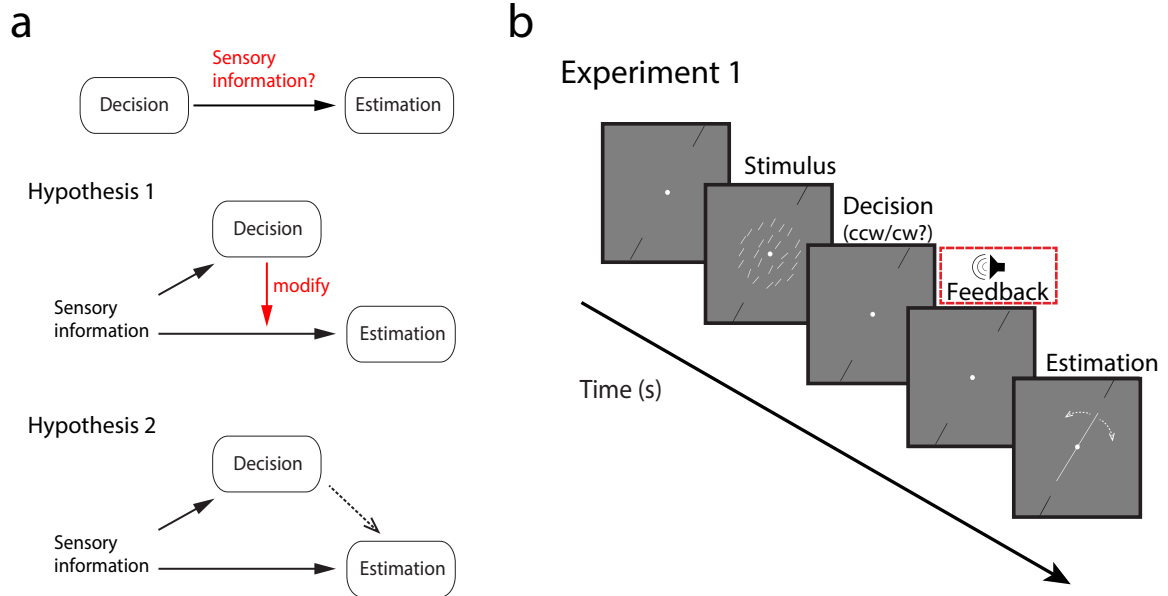
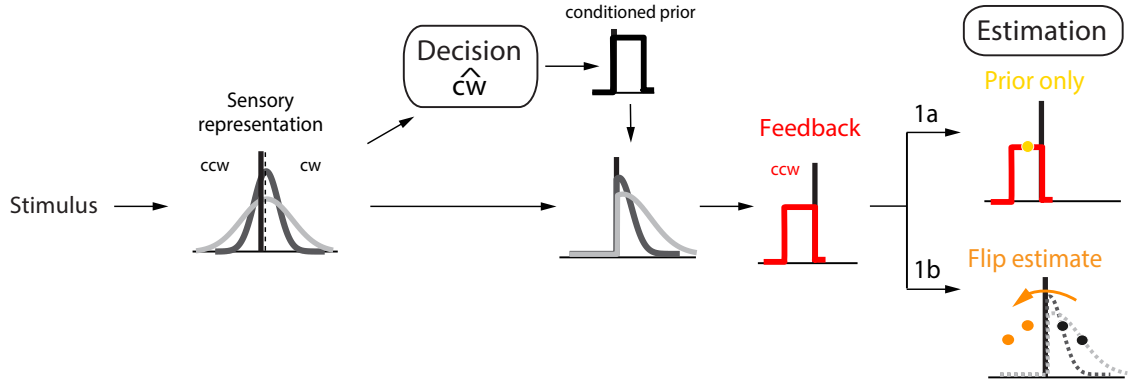


Figure 3.1: *Alternative hypotheses of post-decision sensory representation.* (a) Hypothesis 1: After the decision task, subjects modify the sensory representation in working memory and use this modified representation to make the second inference. Hypothesis 2: Subjects maintain the full sensory representation in working memory and make both decision and estimation based on this representation. (b) Experiment 1: Stimulus is an array of line segments whose orientations are sampled from a Gaussian around a stimulus orientation. The stimulus orientation is sampled uniformly within ± 21 deg of the reference orientation. Subjects first reported whether the stimulus orientation was clockwise or counter-clockwise of a decision boundary indicated by two black lines. After 300 ms, they received an auditory feedback to indicate whether their decision was correct or incorrect. Then they had to take into account the feedback to estimate the stimulus orientation by adjusting a probe line.

server combines the updated prior with the full sensory representation to estimate stimulus orientation. In the second strategy (2b-Resample), because the observer knows that the original sensory measurement is on the incorrect side, it tries to reconstruct a new sensory measurement that falls on the correct side. More specifically, the observer resamples from a distribution centered on the original sensory measurement until it gets a sample that is on the correct side. Then the new sample and the updated prior are used for the estimation. So the key difference between the two strategies is the sensory sample that is used to make the inference.

a

H.1: Stimulus information is erased



b

H.2: Stimulus information is maintained

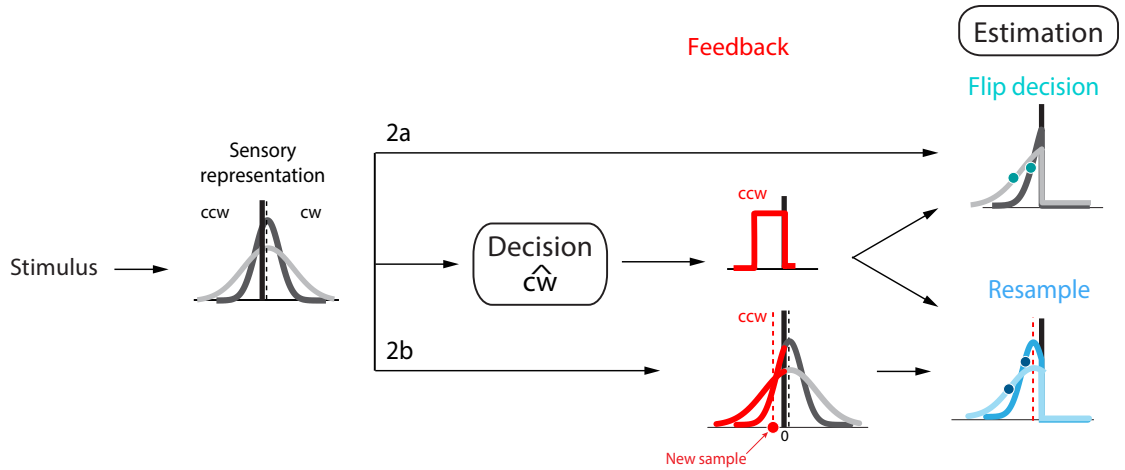


Figure 3.2: *Prediction of alternative hypotheses for incorrect trials.* (a) Hypothesis 1: After making a decision from the sensory representation (illustrated by Gaussian distributions whose width are proportional to stimulus uncertainty), the observer erases the inconsistent sensory information by multiplying the conditioned prior with the sensory representation. Upon receiving feedback indicating that the decision is incorrect, the conditional prior is updated. Strategy 1a assumes that the observer takes the mean of the updated prior as the estimate. For strategy 2a, the observer computes the mean of the modified sensory representation and flipped it to the correct side. (b) Hypothesis 2: The observer preserves the full sensory representation and updated the conditional prior upon receiving the feedback. Strategy 2a posits that the observer combines the full sensory representation and the updated prior to make an estimation. In strategy 2b, the observer constructs a resampling distribution that has the center at the original sensory measurement and has the width proportional to sensory uncertainty. It then draws samples from this distribution until a sample falls on the correct side. This sample is then used to make an estimation.

3.2.2. Some sensory information is still preserved

To test the prediction of the above strategies, I first fit the self-consistent Bayesian observer model (Fig. 3.3a) jointly to all data in the decision task and data for correct trials in the estimation task. The model is the same as the one in Experiment 1, chapter 2. The only difference is that for the current experiment, I only fit the correct trials for the estimation data.

Then I implemented the above strategies using the fit parameters to predict incorrect trials in the estimation task. Note that this is a prediction with no free parameter and the models can predict the full distribution of subjects' estimates in the incorrect trials. A qualitative comparison between the model prediction and the data shows a few noteworthy features (Fig. 3.4). Because the Prior Only model (1a) is independent of sensory noise, it predicts that the estimation distribution is the same across the two stimulus noise levels. The slight difference is due to different proportions of incorrect trials in the decision task. In contrast, the data (combined subject, $n=5$) show that the estimates have larger means and larger variances for higher noise levels. The remaining three models seem to qualitatively capture this pattern in the data. A closer examination of the three models shows that the Flip Estimate (1b) and the Resample (2b) models can better capture the magnitude of subjects' estimates, especially for high noise condition. On the other hand, the Flip Decision (2a) and the Resample (2b) models can better account for the overall trend in subjects' estimates which is roughly constant with a slight increase in magnitude when the stimulus is farther away from the decision boundary. However, the data are too sparse for those stimulus orientations to draw a definitive conclusion.

I used two different metrics to quantitatively compare the models. First, because the models can predict the full estimate distribution, I can compute the likelihood of each model given the data. I also obtained the fit parameters for each individual subject so that I can compare the models at the level of individual subject (Fig. 3.5a). To get a sense of how well the models perform in absolute term, I also consider the oracle model which was derived from the empirical distribution of the data. That serves as the absolute upper bound on the

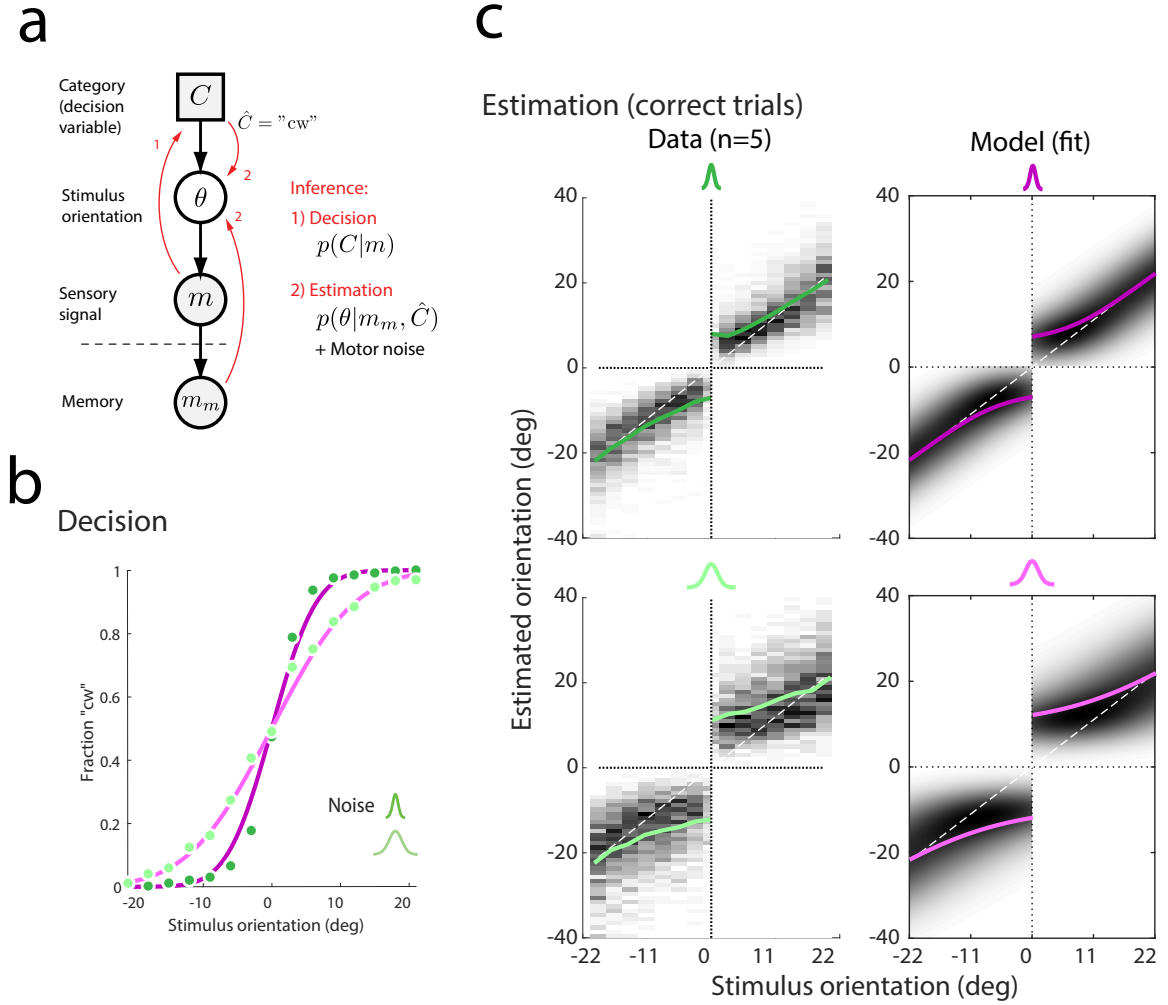


Figure 3.3: *Self-consistent model fit to correct trials.* (a) In each trial, a sensory measurement m is made from the stimulus orientation θ . The observer uses this measurement to make a categorical decision \hat{C} . By the time the observer performs the estimation, the sensory measurement is further degraded in working memory which results in a sample m_m . Given this sample and the conditional prior $p(\theta|\hat{C})$, the observer estimates the stimulus orientation. (b) Model fit to subjects' decision (combined subject, $n=5$). Green circles indicate the data and purple lines indicate the model fit. Darker shade represents low stimulus noise. As expected, the slope of the psychometric function is steeper for lower noise level. (c) Model fit to the estimation data. Subjects show the characteristic biases away from the decision boundary as observed in Chapter 2. The model can capture this bias pattern well. The data and model are shown for correct trials only which explains why there is no estimate in the upper left and lower right regions. Note that the model was jointly fit to the decision and estimation data.

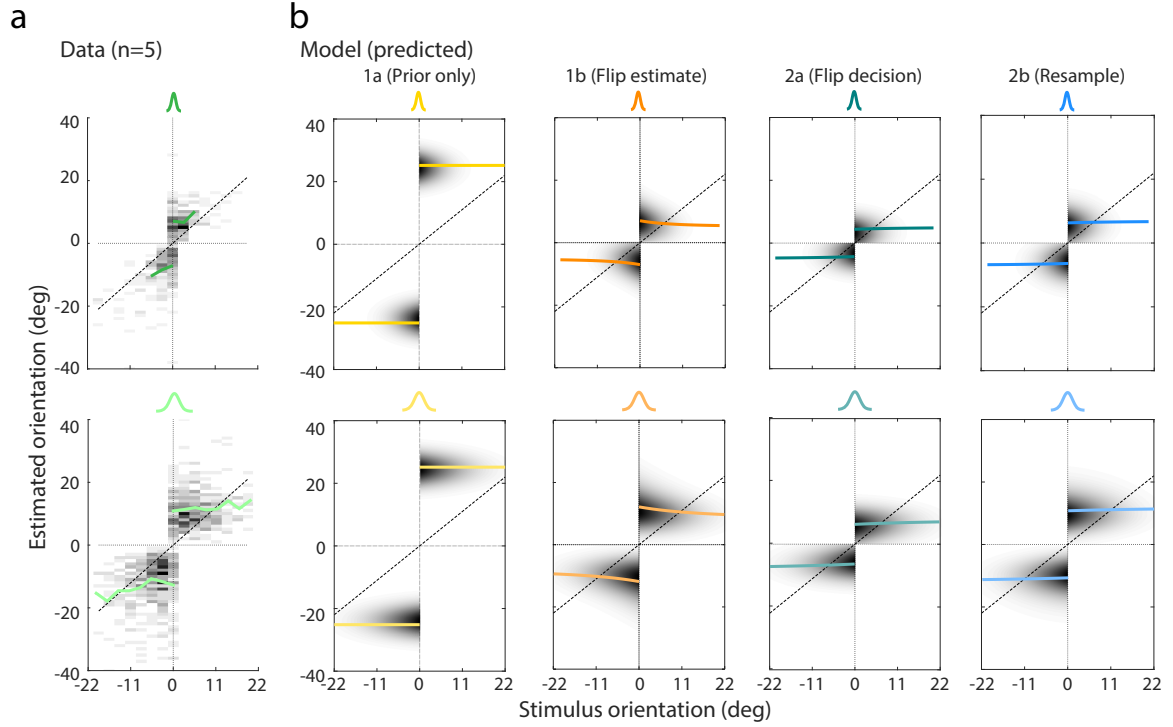


Figure 3.4: *Data and model prediction for incorrect trials.* (a) Distribution of subjects' estimates in incorrect trials. Each row corresponds to a stimulus noise level (upper: low noise). The data is sparse when the stimulus orientation is further away from the decision boundary and the stimulus noise is low because subjects made less incorrect decisions. The mean of subjects' estimate is roughly constant on each side of the decision boundary. Also, the mean and variance of the estimate are larger for higher noise level. (b) The Prior Only model predicts a constant estimate on each side of the boundary. Variation in the model's estimates is only due to motor noise which is the same across all stimulus conditions. Therefore, the estimates' mean and variance are the same across two stimulus noises. In contrast, the remaining three models predict larger mean and variance for higher stimulus noise. However, the Flip Decision model tends to underestimate the magnitude of subjects' estimates, especially in the high noise condition. Moreover, the Flip Estimate model predicts decreasing estimate magnitude when stimulus orientation is further away from 0 which is opposite to the overall trend in the data.

model performance. The result is quite consistent across individual subjects and concordant with the above qualitative comparison between models (Fig. 3.5b). First, the Prior Only model is worse than other models by a large margin. Second, although the Resample model has the highest log likelihood, its likelihood value is not significantly higher than the Flip Estimate and the Flip Decision models. Moreover, bootstrapping the model likelihood also indicates that the Prior Only model is statistically inferior to other models while the

three models Flip Estimate, Flip Decision and Resample are statistically indistinguishable. In the second metric to assess model performance, I computed the mean estimates in each condition for each model and then calculated the correlation and mean squared error (MSE) between model predictions and data. Although the correlations are pretty high and similar across models (0.74-0.82), MSE of the Prior Only model is remarkably higher than the other models (110 vs. 10, 29 and 8). That is clearly observed in the scatter plot of model prediction and data (Fig. 3.5c). Hence the MSE metric also corresponds well with the log likelihood metric and the qualitative analysis.

The above analyses are quite consistent and show that I can rule out the Prior Only model. Because the Prior Only model is the only model that doesn't use any sensory information in the estimation whereas other models do to a certain extent, it suggests subjects still maintain some sensory information when performing the estimation task. However, the analyses cannot clearly dissociate the three models Flip Estimate, Flip Decision and Resample. Therefore, I conducted Experiment 2 to better distinguish those models.

3.3. Further test of alternative models

3.3.1. Imposing asymmetric prior

Experiment 2 was similar to Experiment 1 except for three differences. First, the stimulus noise for low noise condition was slightly higher in Experiment 2. Second, the stimulus range in Experiment 2 was asymmetric around the reference (12 deg on the CCW side and 30 deg on the CW side) whereas the stimulus range in Experiment 1 was symmetric (21 deg on both sides). Third, the subjects were trained to learn this asymmetric stimulus range and were explicitly reminded about the range at the beginning of each trial (Fig. 3.6a). If the subjects could learn this asymmetric prior, the distribution of subjects' estimates will be shifted towards the CW side for correct trials (Fig. 3.6b). Given this learned prior, the models make markedly different predictions for incorrect trials (Fig. 3.6c). Specifically, the Flip Estimate model predicts a shift in the estimation distribution towards the CCW side. In contrast, the Flip Decision and Resample models predict a shift towards the CW side.

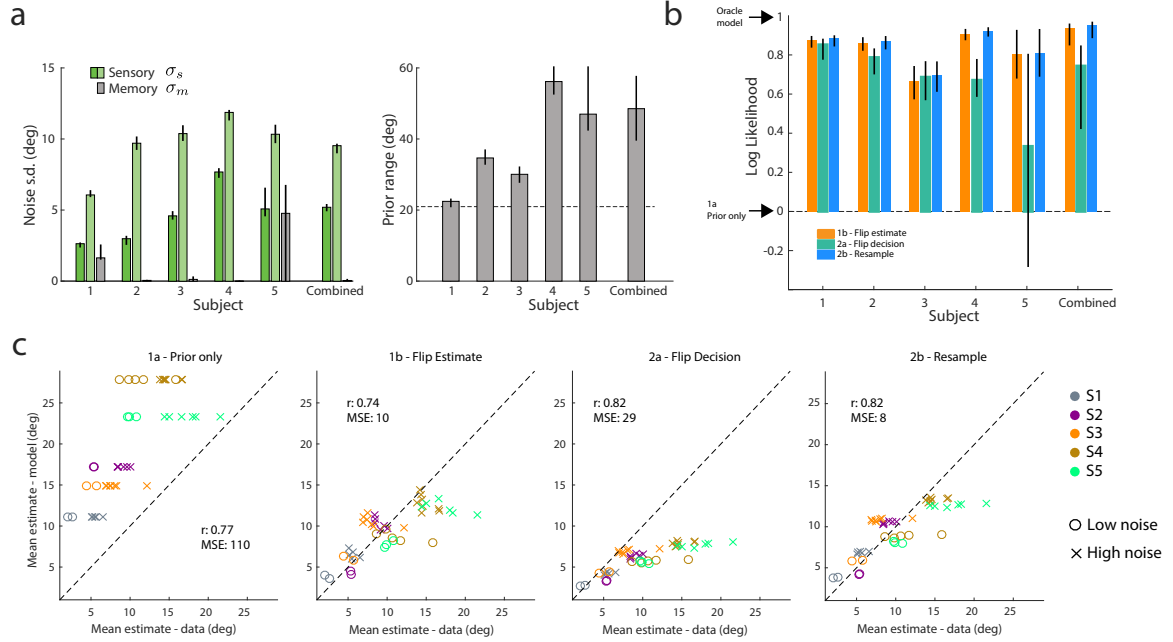


Figure 3.5: *Performance of alternative models in Experiment 1.* (a) The fit noise and prior range show substantial variation across individual subjects. However, it is consistent across subjects that the fit sensory noise increases with stimulus noise and all subjects tend to overestimate the true stimulus range (± 21 deg around the decision boundary). Overall, memory noise is relatively small compared to sensory noise. (b) Normalized log likelihood of the models. Log likelihood values are normalized so that the oracle model (data explaining itself) has a value of 1 and the Prior Only model has a value of 0. It is consistent across all subjects that the Prior Only model performs substantially worse than other models and the Resample model is the best. However, except for subject 4, the performances of Flip Estimate, Flip Decision and Resample models are not statistically different. (c) Scatter plots compare the mean estimates between model and data. Each point corresponds to one experimental condition (stimulus noise and orientation) of one subject. Although all models can roughly explain the trend in data (high correlation), the prediction of Prior Only model is substantially off from the data. Flip Decision model tends to underestimate the magnitude of subjects' estimates. Overall, the performance of Flip Estimate, Flip Decision and Resample models are roughly similar. All errorbars indicate 95% confidence interval computed over 200 bootstrapped samples.

3.3.2. Subjects maintain full sensory representation

To compare the model prediction and the data, I computed the mean estimates on CW and CCW sides for each individual subject and each noise level. I also performed the model fit to correct trials as in Experiment 1 and used the fit parameters to predict the mean estimates in incorrect trials. The fit prior range indicates that all subjects learned the

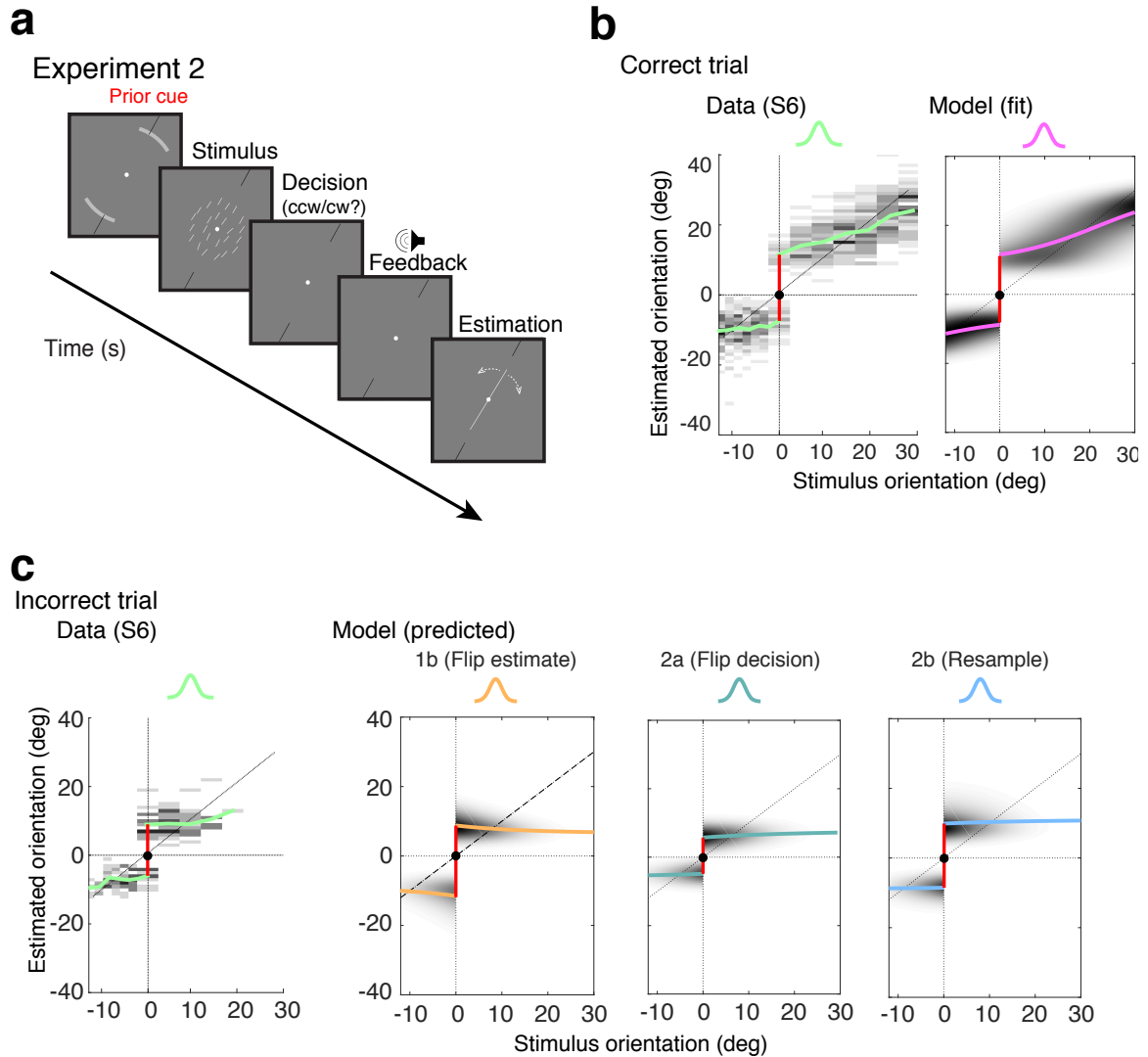


Figure 3.6: *Dissociating the models by imposing an asymmetric prior.* (a) Experiment 2 was similar to Experiment 1 except that the stimulus range was asymmetric around the reference (CCW: 12 deg, CW: 30 deg) and subjects were explicitly reminded about this by two gray arcs shown at the beginning of each trial. (b) Estimate pattern of one representative subject. Estimation data and model fit are shown for correct trials with high stimulus noise. The whole distribution pattern shifts upward to the CW side which suggests that this subject learned the asymmetric stimulus range. (c) Given the learned asymmetric prior, Flip Estimate model predicts a shift in the estimation pattern towards CCW side whereas Flip Decision and Resample models predict an opposite shift towards CW side. Data from this representative subject agree with Flip Decision and Resample models.

asymmetric stimulus range although some learned better than the others (Fig. 3.7a). All except one subject (S7) show a consistent pattern that the mean estimates on CW side are

higher than the mean estimates on CCW side (Fig. 3.7b). In contrast, the Flip Estimate model predicts the opposite pattern for all subjects, that is, the mean estimates are lower on CW compared to CCW side. The Flip Decision and Resample models, however, predict similar pattern as the data.

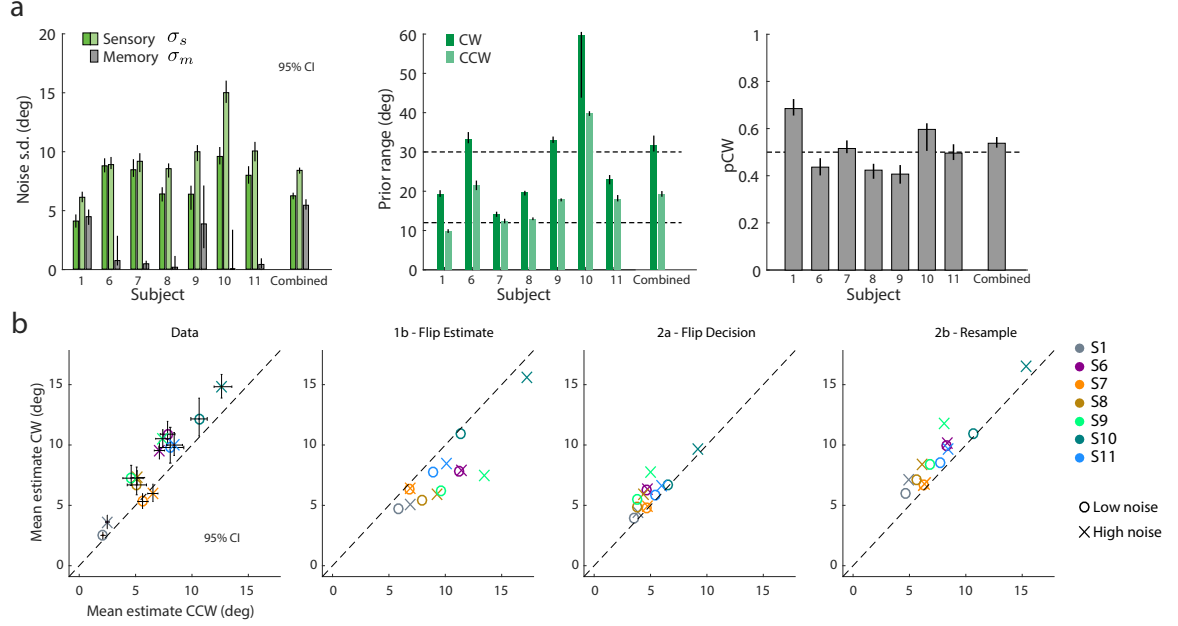


Figure 3.7: *Model prediction and data in Experiment 2.* (a) Fit parameters of individual subjects. Fit sensory and memory noise are comparable to Experiment 1. Sensory noise for the low noise condition is higher because the stimulus noise was increased for low noise condition in Experiment 2. The fit prior range on CW side is higher than on CCW side for all subjects. There is one additional free parameter in Experiment 2 which is subjects' belief about the probability that the stimulus orientation is on CW side (this parameter was fixed at 0.5 in Experiment 1). (b) Comparison of mean estimates between CW and CCW sides for all subjects. Data show that mean estimates are larger on CW compared to CCW side for all except one subject (S7). That pattern is in agreement with Flip Decision and Resample models but is opposite to the Flip Estimate model. Subject S7 seems not able to learn the asymmetric prior (the fit prior ranges are almost identical for CW and CCW sides). Therefore, his data are not informative in distinguishing the models. All errorbars indicate 95% confidence interval computed over 200 bootstrapped samples.

I also computed the models' performance using the metrics as in Experiment 1. The likelihood metric indicates that for 4 out of 7 subjects, both Flip Decision and Resample models outperform Flip Estimate model and for the other 3 subjects, only Resample model outperforms Flip Estimate (Fig. 3.8a). The correlation and MSE metrics are also in agree-

ment with the likelihood metric (Fig. 3.8b). The correlation for Flip Estimate model (0.66) is considerably lower than the correlation for Flip Decision and Resample models (both are 0.92). The MSE for Flip Estimate (7.9) is slightly higher than Flip Decision (7.5) and is substantially higher than Resample model (2.5).

All of the above analyses consistently show that Flip Estimate model is inferior to Flip Decision and Resample models. Because Flip Decision and Resample models rely on the full sensory representation whereas Flip Estimate only relies on the modified sensory representation, the results suggest that subjects preserve the full sensory representation after the categorical decision. Another consistent finding from both Experiments 1 and 2 is that Resample model can best explain the data which suggests an interesting form of self-consistency: subjects reconstruct the sensory evidence to make it consistent with the new information (i.e. the feedback).

3.4. Discussion

Across two experiments, I tested how subjects maintain sensory representation in a sequence of tasks by giving feedback in the first task. Comparisons of alternative models suggest that subjects preserve the full sensory representation after the decision task and use that to make the inference in the subsequent estimation task. Certainly, there are other strategies subjects may use to make an estimation with the modified/intact sensory representation and we cannot exhaustively test all of them. However there is a consistent trend across the two experiments: the more an observer model uses sensory information, the better it can predict the data. Specifically, the worst model is always Prior Only model which uses no sensory information and the best model is always Resample model which relies heavily on the full sensory representation. That points to the same direction of hypothesis 2 that subjects maintain the full sensory representation.

The key difference between the experiments in the current study and previous studies using the same sequential tasks (e.g. Experiment 1, chapter 2) is that subjects were given feedback in the decision task and they knew about that beforehand. So some may argue that subjects' inference will be different compared to the scenario when no feedback is

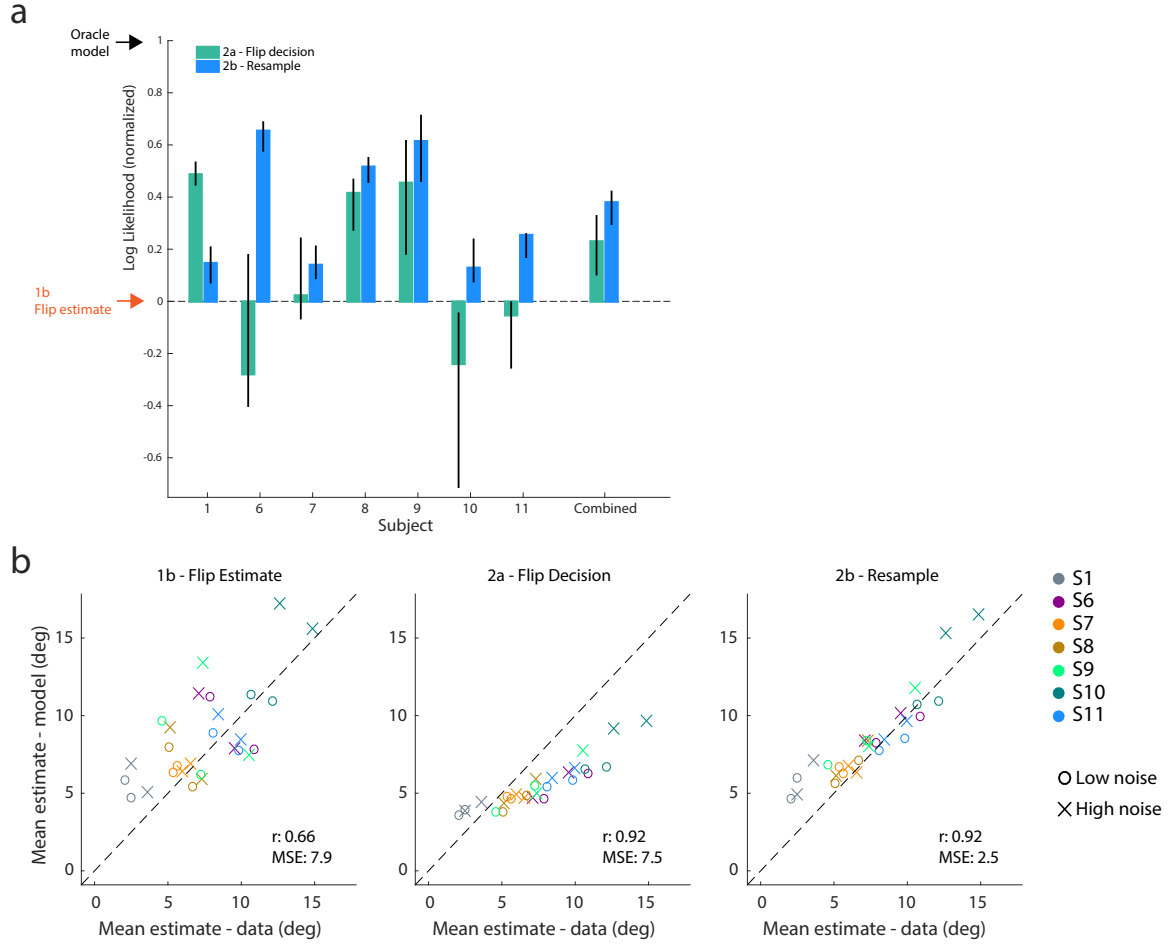


Figure 3.8: *Performance of alternative models in Experiment 2.* (a) Normalized log likelihood of the models. The log likelihood values are normalized such that the oracle model is 1 and the Flip Estimate model is 0. Note that we don't consider the Prior Only model here because it was ruled out in Experiment 1. Hence the lower bound is Flip Estimate instead of Prior Only model. Out of 7 subjects, Resample model is the best for 6 subjects and Flip Decision is the best for 1 subjects. In general, both the Flip Decision and Resample models are better than the Flip Estimate model. (b) Scatter plot of model prediction vs. data. The prediction of Flip Estimate model is worse than Flip Decision and Resample models in terms of correlation and MSE with the data. The prediction of Flip Decision model shows similar underestimation as in Experiment 1. The results are consistent with Experiment 1, that is, Resample model can best explain the data in all metrics. All errorbars indicate 95% confidence interval computed over 200 bootstrapped samples.

given. Thus it may limit the generality of our findings. To solve that issue, we can have an experimental design such that the feedback is only given in one catch trial. Obviously, the limit of that design is there are too few useful trials. Nevertheless, given that the data

patterns in correct trials are similar to the no-feedback scenario and the self-consistent model can quantitatively account well for the data, it is likely that subjects' inference is very similar across the two situations.

It is interesting to see how the sensory representation is maintained in other sequential task settings such as when subjects make a preliminary decision before making a final judgment (Bronfman et al., 2015) or when subjects perform a decision and report confidence on that decision (Pleskac and Busemeyer, 2010; Peters et al., 2017; Yu et al., 2015a). Although these studies found various kinds of sequential dependency between judgments, it is not straightforward whether those effects can be readily explained by the self-consistent Bayesian framework (see section 5.2 for a detailed discussion). Therefore, it is premature to make specific speculation about the mechanism involving the sensory representation in those scenarios.

The results in this chapter may have important implication in more practical situations such as the reliability of eyewitness identification. Several lines of evidence have suggested that human memory is highly malleable and vulnerable to distortion (Wixted et al., 2015; Loftus, 2005) including the self-consistent behavior similar to the choice-induced biases (Loftus, 1975). However, it is unclear whether the representation of original memory (the analog of sensory representation) was modified or not (Loftus and Loftus, 1980). Note that the literature on false memory is mostly about long-term memory whereas my study focuses more on short-term memory. Therefore, there may be significant differences between the two cases.

3.5. Methods

Experimental procedure

Eleven subjects with normal or corrected-to-normal vision (3 males, 8 females; one non-naïve) participated in the experiments. All of the subjects provided informed consent. The experiments were approved by the Institutional Review Board of the University of Pennsylvania under protocol #819634.

General procedure: During the experiments, subjects sat in a darkened room and viewed

the stimuli on a special purpose computer monitor (VIEWPixx3D, refresh rate of 120 Hz and resolution of 1920 x 1080 pixels) at a distance of 83.5 cm (Experiment 1) or 91 cm (Experiment 2). We enforced the viewing distance with a chin rest. All experiments were run in Matlab (Mathworks, Inc.) using the MGL toolbox (<http://justingardner.net/mgl>) and Psychophysics Toolbox (Brainard and Vision, 1997). We used an Apple Mac Pro computer with Quad-Core Intel Xeon 2.93 GHz, 8GB RAM to run Matlab code. Subjects were asked to fixate at the fixation dot whenever it appeared on the screen. Before the main experiments, we trained each subject in 2-3 sessions to familiarize themselves with the task and to have them learned an asymmetric stimulus range (Experiment 2). After that, each subject either completed 2100 trials in 3-5 sessions for Experiment 1 or completed 1820 trials in 3-4 sessions for Experiment 2. Each session lasted approximately 50 minutes. In total, there are 70 trials for each stimulus condition which includes 15 (Experiment 1) or 13 (Experiment 2) stimulus orientations and two noise levels. Subjects used a gamepad (Sony PS4 Dualshock) to give their answers. In the discrimination task, they indicated their decision by pressing a trigger button (left for 'ccw', right for 'cw'). In the estimation task, they reported the perceived stimulus orientation by adjusting a reference line (length: 5 deg) with the analog joystick of the gamepad and then press a button to confirm. Screen background luminance was 40 cd/m² and mean stimulus luminance was 49 cd/m².

Experiment 1: Five subjects (S1-5) participated in Experiment 1. At the beginning of each trial, subjects were presented with a fixation dot (diameter: 0.3°) and a decision boundary indicated by two black lines (length: 3°, distance from fixation: 3.5°). We randomly sampled the orientation of the decision boundary from 0 to 180 (the full circle) in each trial. After 1 s, we presented an orientation stimulus consisting of white line segments (length: 0.6°) that are arranged on two circular arrays centered at the fixation: the outer circle (diameter: 3.8°) has 16 line segments and the inner circle (diameter: 1.8°) has 8 line segments. Small random jitters (from -0.15° to 0.15°) were independently added to the x-y coordinates of each line segment. We sampled the orientation of each line segment from a Gaussian distribution whose mean (stimulus orientation) ranges from -21 (CCW)

to 21 (CW) deg in steps of 3 and standard deviation (stimulus noise) are 3 and 18 deg. The stimulus disappeared after 500 ms and subjects were instructed to indicate whether the stimulus orientation was clockwise or counter-clockwise of the decision boundary. If subjects' response was within 4 seconds, a feedback tone (100 % valid) was briefly played to specify whether subjects' decision was correct (high-pitched) or incorrect (low-pitched). Subjects then went on to report their perceived stimulus orientation. They were instructed explicitly that the feedback was always valid and that they had to take into account the feedback when performing the estimation task. If their response in the decision task was more than 4 seconds, the current trial was skipped and was moved to the end of the trial queue. At the end of each trial, a blank screen (mean luminance) was displayed with a duration randomly chosen from 300 ms to 600 ms.

Experiment 2: Seven subjects (S1, S6-11) participated in Experiment 2. The procedure in Experiment 2 was the same as in Experiment 1 except the following differences. The stimulus range was from -12 to 0 in steps of 2 for the CCW side and from 5 to 30 in steps of 5 for the CW side. We used different step sizes for the two sides to keep the number of stimulus orientation on the CCW and CW sides the same while having a wider range on the CW side. Because it is essential that subjects learned this asymmetric stimulus range, we displayed the stimulus range at the beginning of each trial for both training and main experiments. We also set the standard deviation of low-noise condition to 6 deg which is higher than in Experiment 1 (3 deg) because it increases the discriminability between models.

Training - simple motor task: In the first training, subjects were first presented with a fixation dot. For those participating in Experiment 2, we also displayed a decision boundary (like in the main experiment) and a gray arc indicating the asymmetric stimulus range. After the decision boundary and the gray arc (for Experiment 2 only) went away, a single white line (like the reference line subjects used to do the estimation in main experiments) was shown for 500 ms. Subjects then had to reproduce the presented white line by adjusting a reference line with the analog joystick. After they confirmed by pressing a button,

the original stimulus was displayed again in green on top of subjects' estimates. The decision boundary was uniformly sampled around the circle and the stimulus orientation was uniformly sampled in the same range as in Experiment 1 and 2. Subjects completed 450 trials (Experiment 1) or 325 trials (Experiment 2). We computed the standard deviation of subjects' estimates and used that as a measure of individual motor noise in the model fit and model prediction described below.

Training - estimation task: The estimation training is identical to the motor training except that we used the array stimulus in the main experiment instead of the single line stimulus. Subjects completed 100 trials (Experiment 1) or 150 trials (Experiment 2) for this training.

Training - decision task: The stimulus was identical to estimation training but in this training subjects performed a decision task (CCW/CW) and was given feedback like in the main experiments. Subjects completed 900 trials (Experiment 1) or 200 trials (Experiment 2).

Self-consistent Bayesian observer model

The basic model was similar to formal description in chapter 2, Methods. The difference is in the incorrect trials. For completeness, I describe the model for both correct and incorrect trials.

Decision task

We assume the observer forms a hierarchical model about the task as illustrated in Fig. 3a. Let $C = \{'cw', 'ccw'\}$ be the decision variable indicating whether the stimulus orientation is cw or ccw relative to the decision boundary and θ be the stimulus orientation. In each trial, the observer makes a sensory measurement m from the stimulus orientation θ . We assume the observer knows the sensory uncertainty which is represented by the likelihood function $p(m|\theta)$. The observer can then compute the likelihood over the decision variable as follows:

$$p(m|C) = \int_{-\pi}^{\pi} p(m|\theta)p(\theta|C)d\theta . \quad (3.1)$$

where $p(\theta|C)$ represents the observer's prior expectation about the stimulus orientation on

the two sides of the decision boundary. The observer then computes the posterior over the decision variable by combining the likelihood and the prior $p(C)$:

$$p(C|m) = \frac{p(m|C)p(C)}{p(m)} \quad (3.2)$$

Given a symmetric loss function, the observer chooses the category with higher posterior value (MAP estimate):

$$\hat{C}(m) = \operatorname{argmax}_{C \in \{',_{cw}', ',_{ccw}'\}} p(C|m) . \quad (3.3)$$

We can obtain the model prediction of the psychometric function by marginalizing over the sensory measurement distribution:

$$p(\hat{C}|\theta) = \int p(\hat{C}|m)p(m|\theta)dm . \quad (3.4)$$

Estimation task

After the categorical decision, the observer performs inference for the estimation task. Because it took some time for the observer to finish the decision task, we assume that the sensory measurement m is degraded in working memory and produces a memory sample m_m at the time the observer performs the estimation task. This process is formalized by a draw from the distribution $p(m_m|m)$. The memory-corrupted sensory representation can then be described with the likelihood:

$$p(m_m|\theta) = \int p(m_m|m)p(m|\theta)dm . \quad (3.5)$$

Importantly, the observer takes its categorical decision as a new piece of information. As a result, the prior is updated using the preceding decision which results in a new conditioned prior $p(\theta|\hat{C})$. Here there are 2 ways the observer may store the sensory representation in working memory.

Hypothesis 1: the sensory representation is modified

Because the observer trusts its decision as a fact, it may directly modify the sensory representation by combining the likelihood and the conditioned prior:

$$p_{new}(m_m|\theta) = p(m_m|\theta)p(\theta|\hat{C}) \quad (3.6)$$

Effectively, the observer erases the part of sensory representation that is inconsistent with the decision. If the feedback then indicates that the observer's decision was correct, it then computes the posterior from this modified sensory representation:

$$p(\theta|m_m, \hat{C}) \propto p_{new}(m_m|\theta) \quad (3.7)$$

Assuming L_2 loss function, the estimate is the mean of the posterior:

$$\hat{\theta}(m_m, \hat{C}) = \int_{-\pi}^{\pi} \theta p(\theta|m_m, \hat{C}) d\theta . \quad (3.8)$$

However, if the feedback indicates that the decision was incorrect, the observer cannot perform the computation as above because the sensory representation on the correct side (i.e. the part that is inconsistent with the decision) was erased. We consider two strategies the observer may use to do the estimation.

1a-Prior only: The observer totally disregards the sensory representation and performs inference only on the prior. More specifically, it uses the feedback to update the conditioned prior to obtain $p(\theta|\neg\hat{C})$. Because the observer considers the sensory information as uninformative, the posterior is the same as the prior. Using L_2 loss function as before, the estimate is the mean of the new prior.

1a-Flip Decision: The observer pretends that the decision was correct and performs the inference as in the correct trial to obtain the point estimate $\hat{\theta}(m_m, \hat{C})$. Then because this estimate is on the incorrect side according to the feedback, the observer flips this estimate to the other side which is equivalent to taking $-\hat{\theta}(m_m, \hat{C})$ as the final estimate.

Hypothesis 2: the sensory representation is preserved

The observer preserves the full sensory representation in working memory. In the correct trials, the observer computes the posterior from the likelihood and the conditioned prior:

$$p(\theta|m_m, \hat{C}) = \frac{p(m_m|\theta)p(\theta|\hat{C})}{p(m|\hat{C})} \quad (3.9)$$

The estimate is then the mean of the posterior under L_2 loss function as in hypothesis 1. Computationally, the estimates for correct trials are the same for both hypotheses. The difference only arises in the incorrect trials. Here we consider two strategies to make an estimation in incorrect trials for hypothesis 2.

2a - Flip Decision: Because the categorical decision is incorrect, the observer flips the decision bit and uses that to update the conditioned prior $p(\theta|\neg\hat{C})$. Then it combines this updated prior with the likelihood to make an estimation as in the correct trials.

2b - Resample: The observer also updates the conditioned prior based on feedback as in Flip Decision model. However, because the feedback implies that the original sensory measurement was incorrect, the observer tries to reconstruct a new sensory sample that is consistent with the feedback. In particular, the observer keeps resampling from the distribution $p(m_r|m)$ until it gets a sample that is on the correct side. Then the likelihood constructed around this new sample $p(m_r|\theta)$ is combined with the conditioned prior to make the inference as in the correct trials.

To obtain the predictive distribution, we marginalize over the distribution of memory sample (or the resampled sample for Resample model):

$$p(\hat{\theta}|\theta, \hat{C}) = \int p(\hat{\theta}|m_m, \hat{C})p(m_m|\theta)dm_m . \quad (3.10)$$

and then marginalize over the distribution of decision outcome (the psychometric function):

$$p(\hat{\theta}|\theta) = \sum_{\hat{C}} p(\hat{\theta}|\theta, \hat{C})p(\hat{C}|\theta) . \quad (3.11)$$

To account for subjects' motor noise, we convolve this predictive distribution with a

noise kernel that is specified from the motor noise training experiment.

Model specification

The parametric features of the model components are described below:

- The category prior $p(C)$ is set to 0.5 for Experiment 1 because both categories are equally likely. However, we leave this as a free parameter for Experiment 2 because the asymmetric stimulus range may induce a non-uniform category prior.
- The prior distribution over stimulus orientation $p(\theta|C)$ are either symmetric (Experiment 1) or asymmetric (Experiment 2) around the decision boundary. More specifically, $p(\theta|C)$ is constant from the decision boundary up to a certain range on each side and then monotonically decreases to zero with a cosine roll-off. So the prior distribution is characterized by two parameters: the prior range α which is the distance from the decision boundary to half the magnitude of the uniform range and a smooth factor β which indicates how smooth the roll-off is.
- The measurement distribution $p(m|\theta)$ is a Gaussian that is centered on the stimulus orientation θ and has the standard deviation σ_s proportional to the sensory uncertainty. We assume the sensory uncertainty of each subject only depends on the stimulus noise and is the same across other experimental conditions.
- We assume the original sensory measurement m is corrupted by memory noise and becomes m_m . This is modeled by drawing from a Gaussian $p(m_m|m)$ centered on the original sensory measurement m . The width of this distribution σ_m is assumed to be different across subjects but is the same for each subject across all experimental conditions.
- For the Resample model, the resampling distribution $p(m_r|m_m)$ is centered on the memory sample m_m and has the width $\sqrt{\sigma_s^2 + \sigma_m^2}$.
- The motor noise is modeled as a Gaussian with standard deviation σ_0 that is extracted for each individual subject from the motor noise training experiment.

Model fit

We performed a joint fit to all trials of the decision task and correct trials of the estimation task by maximizing the likelihood of model given the data:

$$p(D|\rho) = \prod_{i=1}^n P(D_i|\rho) = \prod_{i=1}^n P(\hat{C}_i|\rho)p(\hat{\theta}_i|\rho) \quad (3.12)$$

where D is the data, ρ is the parameters of the model, \hat{C}_i is subjects' decision, $\hat{\theta}_i$ is subjects' orientation estimate, i is the trial index and n is the number of trials.

We used Nelder-Mead simplex optimization algorithm to minimize the negative log likelihood $-\log(p(D|\rho))$. We ran the optimization routine thirty times starting at randomized initial parameter values to find the best parameter set.

The fit parameters were then used to make prediction in the incorrect trials by implementing different strategies (see above).

Chapter 4

Self-consistency in number judgment

4.1. Introduction

Previous chapters and other studies (Jazayeri and Movshon, 2007; Zamboni et al., 2016) have shown the post-decision biases in sequential judgments for several low-level visual stimuli such as motion direction and orientation. The biases were well explained by the self-consistent Bayesian model across several experimental settings. Because the bias effect bears a close resemblance to cognitive biases in value-based decision-making (see section 1.2), self-consistent inference is likely to hold for other high-level stimuli as well. In this chapter, I examine whether the bias effect also occurs for a high-level cognitive variable in a similar decision situation.

Human perception of numbers has been shown to be a high-level variable that is independent of many associated low-level features (Anobile et al., 2016; Cicchini et al., 2016; Harvey et al., 2013). Several perceptual effects found in low-level variables also hold for number perception (e.g. adaptation - Burr and Ross (2008), Arrighi et al. (2014) and saccadic compression - Burr et al. (2010)). Furthermore, it has been shown that subjects' sensitivity to number stimulus was reduced after making an intermediate categorical decision (Bronfman et al., 2015). Hence it is interesting to see whether the post-decision bias found in previous chapters also hold for number perception.

Across three experiments using both non-symbolic (dot array) and symbolic (Arabic numeral) forms of number, I found a robust bimodal pattern in subjects' estimates which is similar to that shown in the low-level stimuli. Moreover, the self-consistent Bayesian model can provide a good quantitative fit and prediction of the data. The findings suggest self-consistent behavior is quite general across several types of stimuli.

4.2. Dot array stimulus

4.2.1. Choice-induced bias in perceived number of dots

The first aim is to replicate the sequential bias effect in previous chapters using a non-symbolic form of number stimuli. In Experiment 1, subjects made a sequence of judgments on a dot stimulus (Fig. 4.1). The stimulus consists of an array of dots and was presented for 20 ms or 1000 ms. Several low-level features of the dots (size, location, etc.) were randomized across trials to minimize their effects on number perception (see Methods - Experiment 1 for a detailed description). After the stimulus disappeared, subjects first reported whether the number of dots was less or greater than 40. Subsequently, they estimated the exact number of dots by adjusting a probe line.

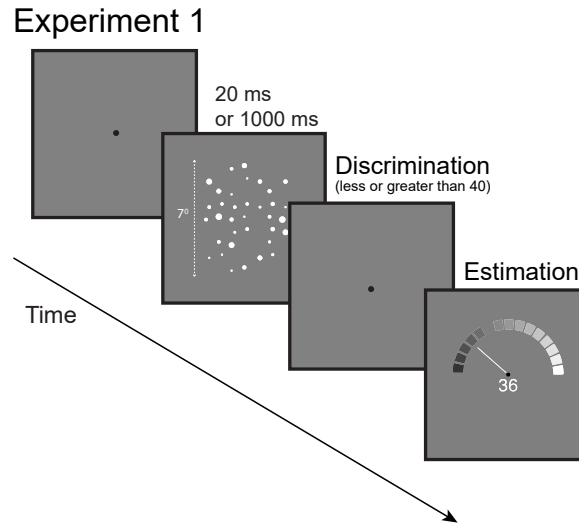


Figure 4.1: *Experiment 1 procedure*. After fixation, subjects were presented with an array of dots. The number of dots was randomly chosen in the range 33 to 47. Presentation time was varied across trials to modulate the task difficulty. Subjects first had to indicate whether the number of dots was less or greater than 40. Then they adjusted a probe line to report the number of dots.

As expected, in decision task subjects performed better when the stimulus duration is longer. Standard deviations of the fit cumulative Gaussian are 3.17 and 4.65 deg for 1000 ms and 20 ms conditions, respectively. There is also a noticeable leftward shift of the psychometric functions (point of subjective equality (PSE): 38.8), suggesting that subjects'

decisions were biased towards greater values (Fig. 4.2a). In the estimation task, subjects exhibited characteristic biases away from the reference at 40 which resembles the bimodal bias patterns using orientation and motion direction stimuli (see e.g. Fig. 2.1). The results suggest that number perception is also subject to similar bias effect as other low-level variables in a sequential task setting.

4.2.2. Explaining the bias with self-consistent model

Given that the dot stimulus results in similar bimodal bias pattern as low-level variables considered before in Chapter 2 and 3, I also expect the self-consistent model to qualitatively capture this pattern. To test whether the model can also quantitatively account for subjects' behavior, I jointly fit the same self-consistent model used in Chapter 2 to both decision and estimation data (see Methods for detailed description of the model). Model fits shows that the basic model can explain the general pattern in the data, yet it fails to capture some nuanced features (Fig. 4.2b). First, it cannot capture the bias in the psychometric function which is expected because I assume a uniform prior over the decision variable, hence forcing the PSE to be at 40. Second, the distribution of estimates predicted by the model is more biased towards the center of the stimulus range than the data. As a consequence, the tail of the distribution is more biased away from the decision boundary than actually found in subjects' estimate. Moreover, subjects' estimates are also generally biased towards higher values which is not captured by the model (Fig. 4.7 in the supplementary material). To test whether the model can be modified to account for those deviations, I consider three variants of the basic model that target three main assumptions: a Gaussian measurement distribution in the linear space, a squared loss function and a uniform prior over decision variable.

Several neuroimaging studies have suggested that the neuronal encoding of number tend to follow a Gaussian shape in the log space (Nieder and Miller, 2003; Nieder and Merten, 2007; Piazza et al., 2004). That is consistent with psychophysics studies showing that human perception of number obeys Weber's law within a certain range (Anobile et al., 2016, 2015). To test how logarithmic encoding of number would change the model prediction,

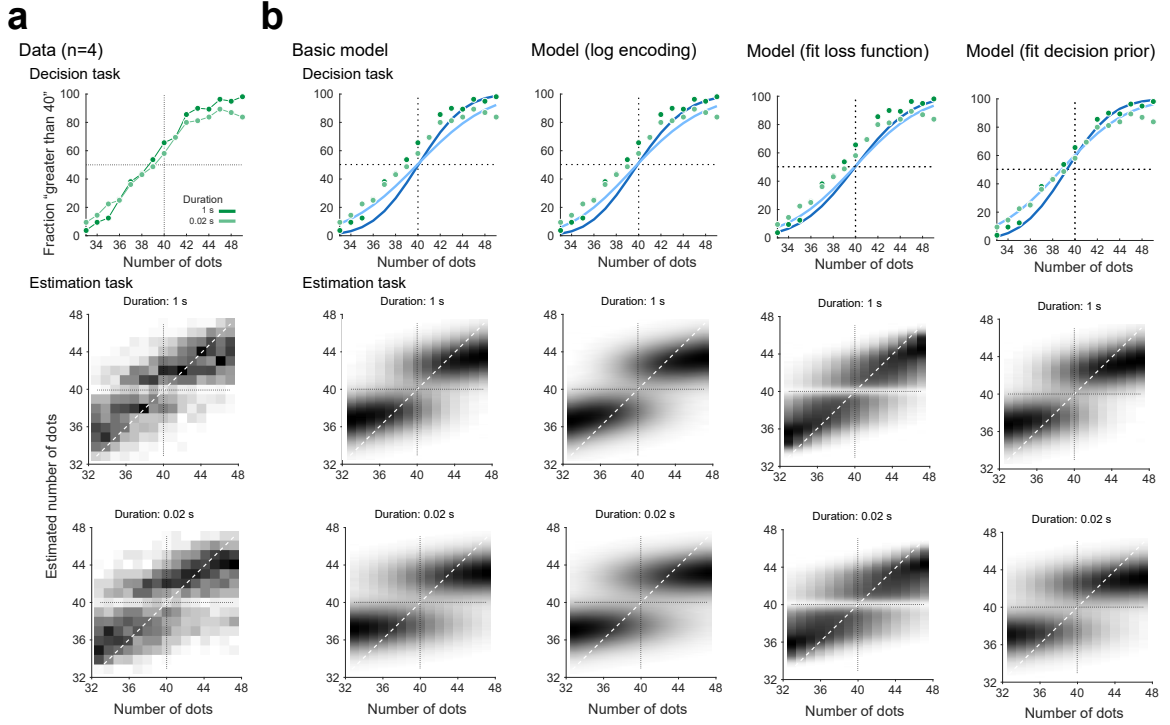


Figure 4.2: *Data and model fit in Experiment 1. (a)* Subjects' data in decision and estimation tasks (combined subject). Psychometric functions indicate that sensory noise is quite similar across the two presentation times. Also, the psychometric functions are shifted to the left, suggesting that subjects' decisions are biased towards 'greater than 40'. Distribution of subjects' estimates shows the characteristic bimodal pattern. The bias magnitudes in subjects' estimates are quite similar for the two presentation times. *(b)* Model fits of different versions of self-consistent model. Basic model uses a Gaussian likelihood in linear space, a squared loss function and a uniform prior over the decision. The other variants change one of those properties. Model with log encoding uses Gaussian likelihood in log space. Model with fit loss function leaves the power of loss function as a free parameter in the fitting. Model with fit prior leaves the prior over decision variable as a free parameter. Although all models can qualitatively account for the bimodal distribution of subjects' estimate, the model with fit loss function provides the best quantitative fit (see Fig. 4.7 for clearer comparison). On the other hand, the model with fit prior best captures the subjects' bias towards "greater than 40" in decision task.

I incorporated it into the basic model by assuming that the measurement distribution is Gaussian in logarithmic space. It effectively results in measurement distributions skewing towards higher values and distribution variance increasing with the stimulus number (Nieder and Miller, 2003). As a consequence, there is an inherent bias towards higher values. Although model fit to subjects' data reflects this bias in the leftward shift of psychometric

function, the bias magnitude is too small to account for the actual bias in the data (Fig. 4.2b). Furthermore, the model still cannot account for the smaller bias in the tail of the estimate distribution. In general, the model fit using logarithmic encoding is very similar to the basic model.

One important component of a Bayesian observer model is the loss function which represents how decision error is penalized. When the stimulus variable is continuous with a well-defined metric space, it is reasonable to assume that large error is penalized more than small error. This is often parameterized by a loss function of power form $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|^n$ in which θ is the stimulus variable, $\hat{\theta}$ is the observer's estimate and n indicates how much to penalize larger error. In the basic model, I assume a squared loss function ($n = 2$) which is a common choice in perception. To test how the model may change when I relax that assumption, I let the power n of the loss function as a free parameter in the model fit. With the fit power close to 0 ($n = 0.0003$), the model can account pretty well for subjects' estimate distribution, especially at the long tail (Fig. 4.2b). It suggests that subjects prioritized getting the exactly correct answer rather than trying to get as close as possible to the correct answer. That may be because in the current experimental task, the stimulus variable comes in discrete form (integer numbers) thus implicitly encouraging subjects to employ the winner-takes-all loss function.

Although the previous modifications to the model could address the estimate distribution, they still fail to account for the bias in subjects' decision towards greater values. A straightforward fix is to change the model's prior over the decision variable. As expected, fitting the decision prior ($p(\text{"greater"}) = 0.57$) helps the model to capture the decision bias by shifting the psychometric function to the left. Interestingly, it can also capture a general small shift in the estimate distribution towards greater values.

In general, the modeling suggests that with appropriate modifications, the self-consistent model can quantitatively account for the subjects' behavior in both decision and estimation tasks.

4.3. Symbolic number stimulus

4.3.1. Probabilistic inference over a sequence of numbers

The post-decision bias found in the previous section using dot stimulus provides supporting evidence for the generality of the sequential effect in high-level stimulus. However, several researchers have been arguing that the sense of number using dot stimulus is still affected by low-level features such as dot density, size, etc. (Gebuis and Reynvoet, 2012; Hurewitz et al., 2006). Therefore, the bias in the perceived number of dots may be confounded by low-level features. To address this issue, I conducted two more experiments using symbolic numbers and designed the task as a form of statistical inference.

In Experiment 2, the stimulus was a sequence of eight Arabic numerals drawn from a truncated discrete Gaussian distribution centered on the stimulus number (Fig. 4.3a, see Methods- Experiment 2 for details). After the stimulus presentation, subjects first decided whether the stimulus number was less or greater than 40 and then provided an estimate of the stimulus number. Because the numerals were presented clearly on the screen, most of the stimulus uncertainty comes from the sampling distribution of the presented numbers. That is in contrast to the dot cloud stimulus whose uncertainty comes mostly from sensory noise. Therefore I have better control over subjects' perceptual uncertainty and can obtain two markedly different uncertainty levels as observed in subjects' decisions (see psychometric functions in Fig. 4.3b). That is also reflected in subjects' estimates: the bias and variance are exceedingly larger for high uncertainty condition (Fig. 4.3b and Fig. 4.3c). The distribution of subjects' estimates reveals a similar bias pattern as in the dot stimulus (Fig. 4.3c). The results suggest that post-decision bias in sequential tasks holds for the number variable, regardless of whether it is symbolic or non-symbolic forms.

4.3.2. Trial-by-trial prediction of observer models

Because the sequence of numerals was shown clearly at a moderate rate, it is reasonable to assume that subjects registered those samples quite accurately and used those samples to make judgments on the stimulus number. Given that assumption, ideal observer models can make predictions on subjects' decision and estimate in each trial (Fig. 4.4a). Note that

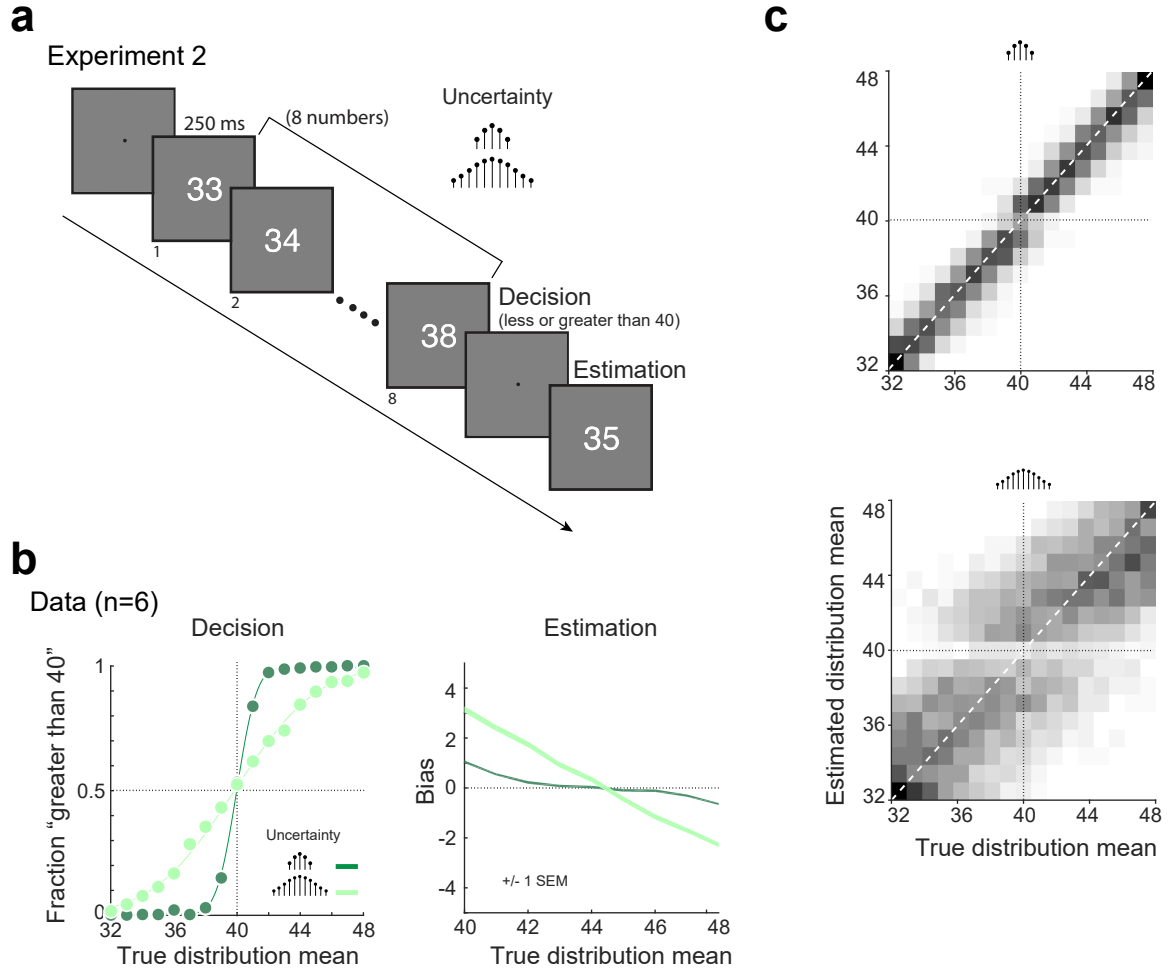


Figure 4.3: *Experiment 2 procedure and result.* (a) Stimulus is a sequence of 8 Arabic numerals drawn from a discrete truncated Gaussian centered on a stimulus number. The width of the Gaussian was modulated to obtain different stimulus uncertainty (standard deviation was 2 or 10). The stimulus number was chosen randomly in each trial within the range 32 to 48. Subjects first indicated whether the stimulus number was less or greater than 40. Then they had to report the stimulus number. (b) Data (combined subject) for decision and estimation tasks. When the stimulus uncertainty increases, the psychometric function is shallower and the estimation bias is larger. Note that the estimate bias is shown for correct trials only. (c) The full distribution of subjects' estimates across two levels of stimulus uncertainty. Contrary to Experiment 1, because subjects' sensory noise is markedly different for two stimulus uncertainties as shown in (b), there is significant differences in estimate distributions across the uncertainty levels. More specifically, the characteristic bimodal pattern is clearly observed in the high uncertainty condition whereas the biases are small for the low uncertainty condition.

for the dot stimulus, most of the uncertainty is from the internal sensory noise. As a result, I don't have access to the sensory samples subjects used to make judgments in each trial, hence the observer models can only predict the overall distribution of decision and estimate.

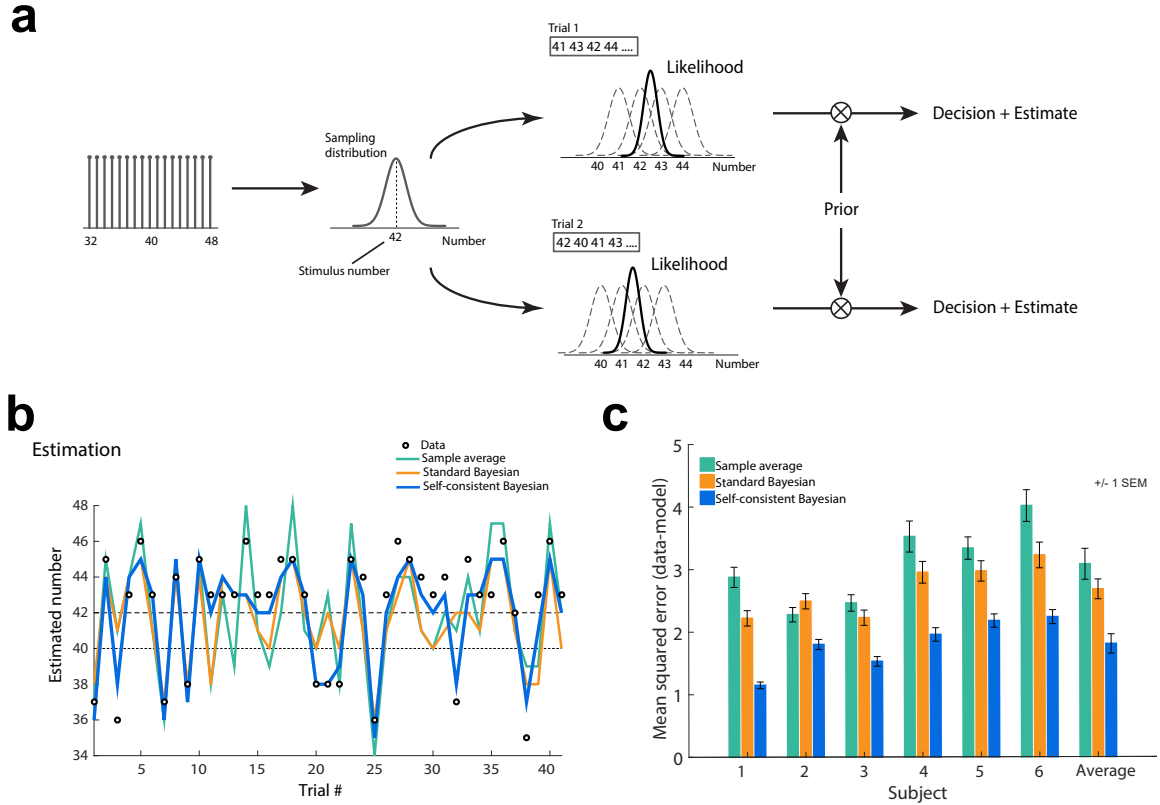


Figure 4.4: *Comparing alternative models.* (a) Trial-by-trial prediction. In each trial, an observer model uses the presented number sequence to make judgments. For example, probabilistic models construct the overall likelihood from the sample numbers and combine it with the prior to make specific predictions in each trial for both decision and estimation tasks. (b) Example data (black circles) from a representative subject (S1) is shown for the condition when the stimulus number is 42 and the stimulus uncertainty is high. The Sample Average model simply takes the mean of presented numbers. The standard Bayesian model makes inference from the prior and likelihood determined by the experimental design. The self-consistent Bayesian model is similar to the standard Bayesian model except that the prior is conditioned on subjects' decision in each trial. Although all models can track subject's estimates relatively well, the self-consistent model is better in some cases (esp. when subjects' decision is different from model prediction as in trial 3). (c) Mean-squared error between model prediction and subject's estimates were computed for each individual subject. The results are quite consistent across all subjects: The self-consistent model performs significantly better than other models and the Sample Average model is the worst. Error bars indicate ± 1 SEM computed over all trials.

Here I consider three alternative models. The Sample Average model assumes that the observer disregards information about stimulus uncertainty and makes judgments using the average of presented numbers. Both standard Bayesian and self-consistent Bayesian model construct the overall likelihood from the sample likelihoods and combine that with the prior to make both decision and estimates. The difference between these two models is the prior. The standard Bayesian model uses the experimental prior which is uniform over the stimulus range. On the other hand, the self-consistent model uses a conditional prior that depends on the subjects' decision in each trial. For example, if subjects' decision is "greater than 40", the conditional prior is uniform over the range [40, 48]. In other words, the self-consistent Bayesian model uses subjects' own decision to predict the estimate on a trial-by-trial basis. The likelihood of two Bayesian models is determined by the stimulus generation process. Consequently, there is no free parameter in the models. Moreover, the prediction of Sample Average model is the same for all subjects. That is also true for the standard Bayesian model but not for the self-consistent Bayesian model because the latter uses subjects' decision to make prediction on the estimate. In general, all models can predict subjects' decisions and estimate relatively well (see Fig. 4.4b for an example). To quantify the predictive power of the models, I computed the mean squared error between model prediction and data across all trials for each subject. The results show a consistent trend across all subjects: self-consistent Bayesian model are better than other models by a significant amount (Fig. 4.4c). Because the key difference between self-consistent model and other models (especially the standard Bayesian model) is its reliance on subjects' decision, the error analysis suggests that subjects' decision play an important role in the estimation process.

4.3.3. Replicating the result with different task instruction

Although the results using symbolic number stimulus provide further support for post-decision bias in number perception, there are two limitations in the experimental design. First, by asking whether the stimulus number was less or greater than 40, subjects may develop a tendency to avoid number 40 in their response, thus artificially resulting in the

bimodal pattern in estimate distribution. Second, the training doesn't explicitly instruct subjects on how the stimulus was generated. Subjects were supposed to learn about the stimulus by practicing the tasks with feedback in each trial. Therefore, I conducted Experiment 3 to eliminate those limitations.

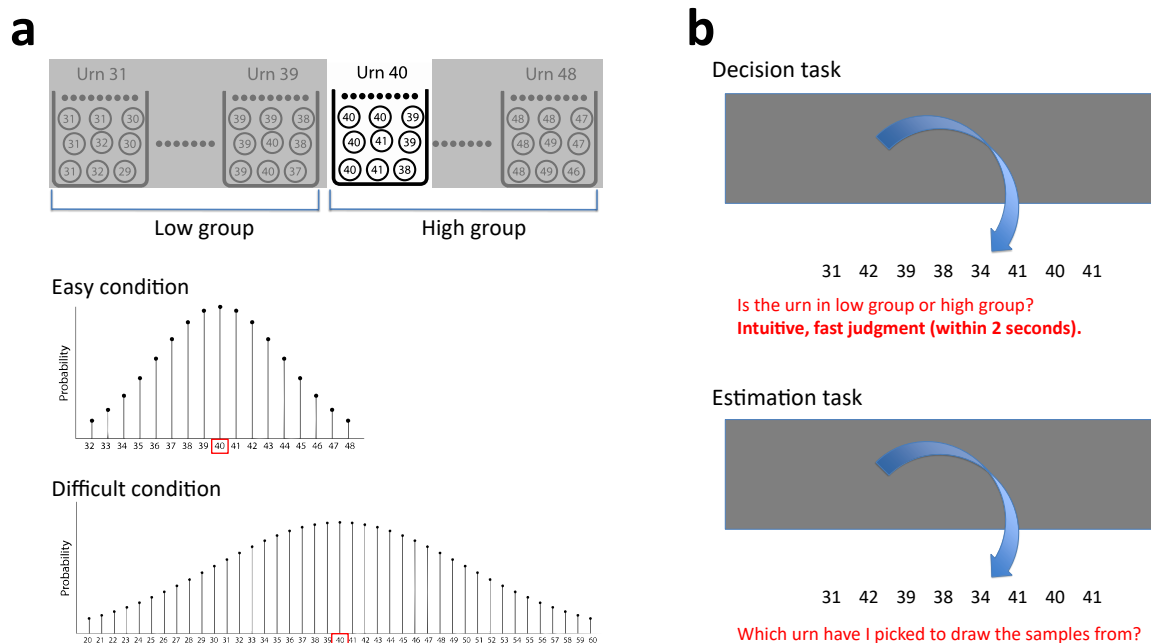


Figure 4.5: *Instruction for Experiment 3.* (a) Stimulus generation process. Eighteen urns contain a large number of balls. The urns are arranged into low (31-39) and high (40-48) groups. The distribution of balls in each urn is centered on the urn label (40 in the illustrative example). The distribution width is determined by the uncertainty condition: narrow for the easy and wide for the difficult condition. (b) Experimental tasks. In each trial, an urn is randomly chosen. Then a sequence of eight balls is drawn from this urn. Subjects had to first decide whether the chosen urn was from low or high group. Subsequently, they had to report the urn label. To prevent subjects from implicitly making the estimate first, it was emphasized that the decision task has to be performed quickly. Note that this is for the training only and the procedure in the main experiment is the same as in Experiment 2 (see Fig. 4.3a)

The design of Experiment 3 was similar to Experiment 2 except for the following differences. First, I split the stimulus numbers into two groups: low group (31-39) and high group (40-48) so that the decision task is whether the stimulus is in low group or high group. That eliminates concern about subjects avoiding the decision boundary at the middle of the stimulus range. Second, I gave subjects an intuitive description of the stimulus

generation and the experimental task in form of an urn example (Fig. 4.5). The instruction emphasizes that the experimental tasks are probabilistic inference over the hidden variable (stimulus number) from a set of observed samples (Arabic numerals). There were two other minor changes in Experiment 3. The stimulus uncertainty was slightly higher in the easy condition and I had separate blocks of easy and difficult trials which were explicitly labeled for subjects. The latter was meant to improve the accuracy of subjects' statistical inference, hence facilitates the comparison between probabilistic observer models (the standard and the self-consistent Bayesian models).

Experimental results confirm what I have found in Experiment 2 with a clear bimodal bias pattern in the distribution of subjects' estimates (Fig. 4.6b). Subjects' judgments in both decision and estimation tasks are also noisier in the easy condition than in Experiment 2, reflecting the increase in stimulus noise in that condition (Fig. 4.6a, b). Prediction errors of alternative models vary widely across subjects, yet follows the same consistent pattern in Experiment 2: the self-consistent model is always the best, followed by the standard Bayesian model. The errors between model prediction and data are also higher compared to Experiment 2. That may be due to the increased stimulus uncertainty and wider stimulus range. Together with the results of previous experiments, the findings suggest subjects did perform probabilistic inference instead of a heuristic strategy (Sample Average) and subjects decision affects the subsequent estimation process.

4.4. Discussion

I have shown that the perceived number of dots was systematically biased away from the decision boundary after having made a categorical decision on the same stimulus variable (Experiment 1). The bias pattern is similar to that found in several low-level stimulus variables. Further experimental tests and model analysis using a sequence of Arabic numerals show similar results in the bias pattern of subjects' estimate (Experiment 2 and 3). The findings suggest that post-decision bias in sequential decision tasks holds for a high-level stimulus. Therefore, self-consistent inference may be a general strategy humans use in many sequential task scenarios.

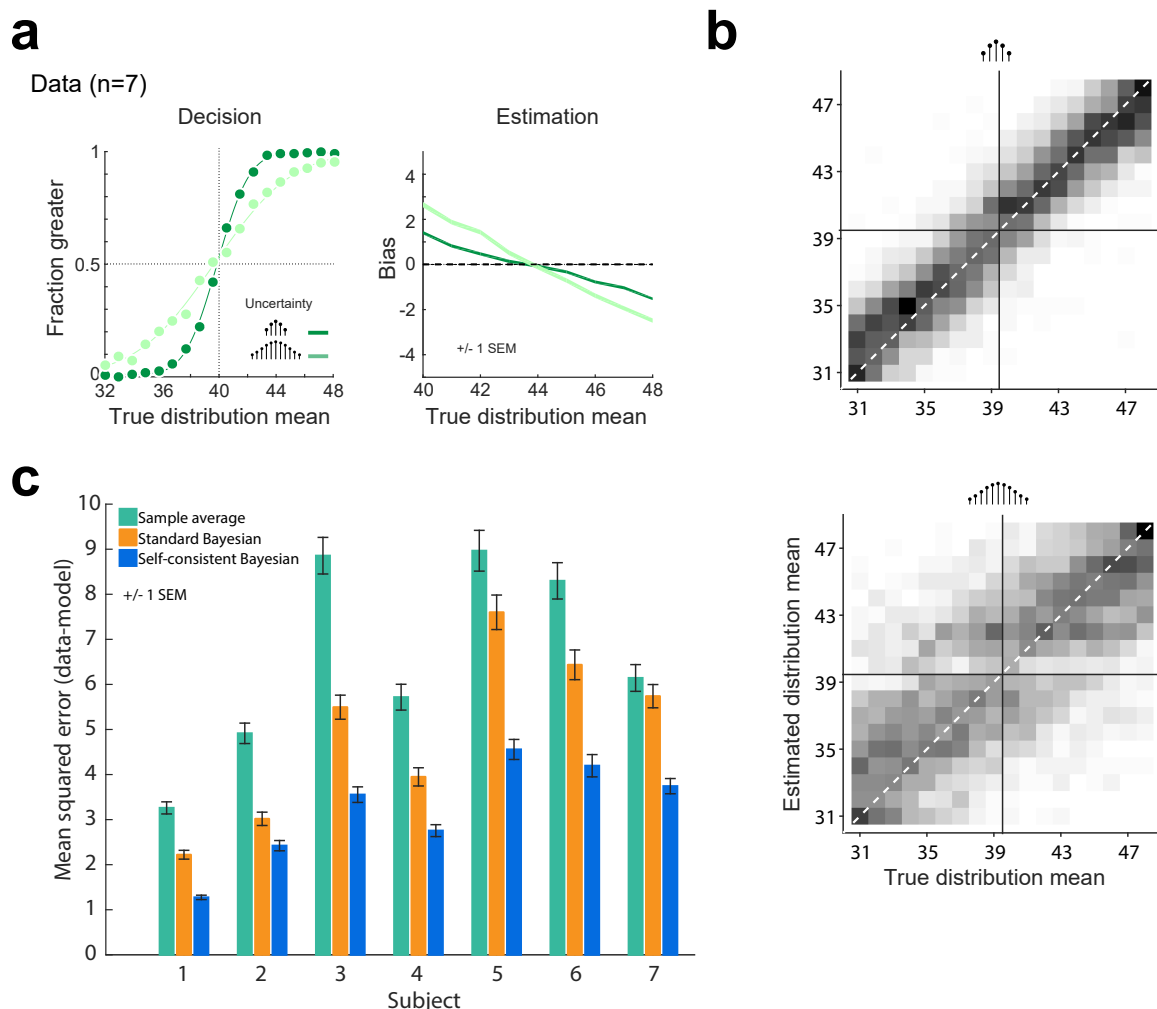


Figure 4.6: *Experiment 3 result.* (a) The results in both decision and estimation tasks are similar to Experiment 2. The psychometric curve in low uncertainty condition is shallower compared to Experiment 2 because a slightly higher uncertainty level was used. (b) The bimodal bias pattern in subjects' estimates are also similar to that in Experiment 2. (c) The prediction errors of alternative models are consistent across all subjects. The self-consistent model is always the best, followed by the standard Bayesian model. Generally, the errors are slightly higher than in Experiment 2.

Model fit to subjects' judgments in dot stimulus experiment suggests that subjects employ a loss function close to L_0 loss (MAP estimator). The similar distribution of subjects' estimate in symbolic number experiments gives further support for that finding. Most studies applying Bayesian model to perception explain experimental data by addressing the prior distribution (Knill and Richards, 1996; Weiss et al., 2002; Girshick et al., 2011; Seriès

and Seitz, 2013) and/or the likelihood function (Wei and Stocker, 2015; Ernst and Banks, 2002; Alais and Burr, 2004; Ernst et al., 2000; Sato and Kording, 2014). The loss function is somewhat neglected in perception studies (but see e.g. Whiteley and Sahani (2008); Landy et al. (2007)). The reason may be partly because the loss function mainly characterizes how an action is rewarded/penalized while perception research is more concerned about subjective experience about the world, not the action that follows a perceptual judgment or the reward that follows the action. Nevertheless, even without a clear reward contingency, the task structure may predispose subjects towards a certain loss function. For example, in the sequential-task experiment with low-level stimuli (e.g. orientation stimulus in chapter 1), subjects provided the estimate on a continuous scale. On the other hand, the number experiments in this chapter required subjects to report their estimate on a discrete scale of integer numbers. As a result, subjects in number experiments may be more inclined towards L_0 loss whereas subjects in orientation experiment are closer to using L_2 loss despite the fact that both stimulus variables are continuous and subjects were instructed to get as close as possible to the correct answer.

In symbolic number experiments (Experiment 2 and 3), all observer models assume that subjects have perfect memory of the presented samples and complete knowledge of the stimulus statistics. In practice, those strong assumptions are rarely satisfied by an average human subject. For example, several studies have shown that human subjects tend to give different weights to sequentially presented stimuli in their judgments. Specifically the weights are higher for early samples (primacy effect) and/or late samples (recency effect) for a wide range of stimuli including words (see Davelaar et al. (2005) for a review), low-level perceptual stimuli (Cheadle et al., 2014; Tsetsos et al., 2012; Drugowitsch et al., 2016; Kiani et al., 2008) and symbolic number (Brezis et al., 2015). In our experiments, there is a slight recency effect in Experiment 3 but no clear pattern was found in Experiment 3 (Fig. 4.8 in the supplementary material). Another possible imperfection is that subjects do not use the exact sampling distribution of the stimulus to make inference. Instead, they may resort to some heuristic to estimate the variance of that distribution (e.g. taking the

variance from the presented samples). Incorporating those constraints into the observer models may improve their prediction. However, it is unlikely that doing so will change the general pattern in the model prediction errors.

Human judgments have been shown to be irrational in a vast amount of high-level cognitive experiments that explicitly stated the probabilistic structure of the task in numerical terms (e.g. choosing between gambles) (Tversky and Kahneman, 1974, 1981; Kahneman, 2011). In contrast, human behaviors are close to optimal inference in a wide range of low-level perceptual and motor tasks. Several researchers have attempted to explain the discrepancy between low-level and high-level experiments (e.g. two decision-making system - Kahneman (2003), Oaksford and Hall (2016) or experience-description gap - Hertwig and Erev (2009), Rakow and Newell (2010)). However, the experimental results in this chapter suggest that in sequential decision settings, the behavioral pattern of human subjects are quite similar regardless of whether the stimulus is low-level (e.g. orientation) or high-level (e.g. a sequence of Arabic numerals). It is actually consistent with more recent studies showing that when the experiments using low-level and high-level stimuli are well controlled and well matched, there is no significant difference in human judgments (Jarvstad et al., 2013, 2012). Therefore, the distinction between cognitive and perceptuo-motor judgments may be more apparent than real.

4.5. Method and supplementary material

Experimental procedure

Fourteen subjects with normal or corrected-to-normal vision (one non-naïve) participated in the experiments. All of the subjects provided informed consent. The experiments were approved by the Institutional Review Board of the University of Pennsylvania under protocol #819634.

General procedure

During the experiments, subjects sat in a darkened room and viewed the stimuli on a special purpose computer monitor (VIEWPixx3D, refresh rate of 120 Hz and resolution

of 1920 x 1080 pixels) at a distance of 87 cm. I enforced the viewing distance with a chin rest. All experiments were run in Matlab (Mathworks, Inc.) using the MGL toolbox (<http://justingardner.net/mgl>) and Psychophysics Toolbox (Brainard and Vision, 1997). I used an Apple Mac Pro computer with Quad-Core Intel Xeon 2.93 GHz, 8GB RAM to run Matlab code. Subjects were asked to fixate at the fixation dot whenever it appeared on the screen. Before the main experiments, subjects had training sessions to familiarize themselves with the experimental tasks so that they could perform the task properly. The training spanned 2-3 sessions (Experiment 1 and 2) or 1 session (Experiment 3). After that, subjects performed the main experiment spreading over several days (1200 trials in 2-3 session for Experiment 1; 1394 trials in 3-5 sessions for Experiment 2; 1476 trials in 4-6 sessions for Experiment 3). Each session lasted approximately 50 minutes. In each experimental condition, there are 40 trials (Experiment 1) or 41 trials (Experiment 2 and 3). Screen background luminance was 40 cd/m² and mean stimulus luminance was 49 cd/m².

Experiment 1

Main experiment: Four subjects (S1-4) participated in Experiment 1. At the beginning of each trial, subjects were presented with a fixation dot (diameter: 0.3°). After 1 s, an array of dots were presented for either 20 ms or 1000 ms. In each trial, I randomly sampled the number of dots from a uniform distribution over the range 33 to 47. The dots' locations were randomly chosen inside a circular aperture of 7 degrees in diameter. The diameter of each dot also randomly varied from 0.07 to 0.35 degrees. The dots' location and size were generated using Matlab code provided in Gebuis and Reynvoet (2011). After the dots disappeared, subjects used a gamepad (Sony PS4 Dualshock) to perform 2 tasks. First, they had to indicate whether the number of dots was less or greater than 40 by pressing a trigger button (left for 'less', right for 'greater'). Then they had to estimate the number of dots by adjusting an analog stick on the gamepad and pressing a confirm button. If their response in the first task was more than 4 seconds, that trial was skipped and added to the end of the trial queue. At the end of each trial, a blank screen (mean luminance) was

displayed with a duration randomly chosen from 300 ms to 600 ms.

Training-Discrimination: The experimental procedure was similar to the main experiment except that (i) subjects only performed the discrimination task (less or greater than 40) and (ii) they were given auditory feedback (100% valid) in every trial (high-pitched beep: correct, low-pitched beep: incorrect). Subjects performed 1200 trials in this training.

Training-Estimation: The experimental procedure was similar to the main experiment except that (i) subjects only performed the estimation task and (ii) the correct number of dots were displayed in green color after their response. Subjects performed 600 trials in this training.

Experiment 2

Main experiment: Six subjects participated in Experiment 2 (S1, S5-S9). After a black fixation dot was displayed for 1 s, a sequence of 8 Arabic numerals was presented one at a time at the center of the screen. Each numeral was presented for 250 ms with a blank screen of 50 ms between consecutive numbers. The numbers were generated as follows. In each trial, a true number was drawn randomly from 32 to 48. Then a sample of 8 numbers was drawn from a truncated discretized normal distribution centered on the true number. To create that modified normal distribution, I first truncated a continuous normal distribution by discarding the tails on the two sides that are more than 2 standard deviations from the mean. Then I discretized the distribution at the integer values and normalized to make it a proper probability distribution. To modulate the stimulus noise, I chose two values for the standard deviations: 2 for low and 10 for high noise. After the last number was presented, a black fixation dot was displayed and subjects had to indicate whether the true number was less or greater than 40 by pressing a key on a computer keyboard ('s': less, 'f': greater). Then subjects reported the true number by adjusting a knob (Griffin PowerMate, wired version) and pressing a key on the keyboard ('Enter') to confirm. Similar to Experiment 1, there was a time-out for the first task so that the trial was skipped if subjects failed to respond within 4 seconds. An inter-trial interval was chosen randomly in the range 500-800 ms.

Training-Discrimination: The experimental procedure was similar to the main experiment except that (i) subjects only performed the discrimination task (less or greater than 40) and (ii) they were given auditory feedback (100% valid) in every trial (high-pitched beep: correct, low-pitched beep: incorrect). Subjects performed 272 trials in this training.

Training-Estimation: The experimental procedure was similar to the main experiment except that (i) subjects only performed the estimation task and (ii) after their response, the true number was shown in green color above subjects' estimate. Subjects performed 272 trials in this training.

Experiment 3

Seven subjects participated in Experiment 3 (S1, S10-S15). The procedure was similar to Experiment 2 except the differences listed below.

Main experiment: The true number in each trial was randomly drawn over the range 31-48. Two standard deviations were used for the sampling distribution: 4 and 10. In Experiment 2, the two stimulus noise conditions are randomly interleaved across all trials. In contrast, Experiment 3 was split into blocks of 10 trials with each block only containing one noise condition. Subjects were explicitly informed about the noise level at the beginning of each block (e.g. "This is an easy block"). In the decision task, subjects indicated whether the true number was in the low group (31-39) or in the high group (40-48). The time-out for the decision task is 2 seconds instead of 4 seconds.

Training: Subjects were given a more detailed and intuitive description of how the stimulus was generated. More specifically, the task was framed in the context of an urn example. Fig.4.5 illustrates this with figures taken from the instruction slides. There are 18 urns each of them is labeled from 31 to 48. The urns are split into 2 groups: low group (urns 31-39) and high group (urns 40-48). In each urn, there is a large number of balls and each ball is marked with a number. The content of each urn is determined by its label. For instance, if the urn's label is 40, the distribution of balls is centered on 40 and the width is proportional to the task difficulty (see Fig.4.5a). Given this setup, subjects' task is as follows (see Fig.4.5b). In each trial, an urn is randomly chosen but subjects do not

know which urn it is. Then the experimenter sequentially draws 8 balls from the urn and shows them to subjects. First, subjects have to guess whether the chosen urn is in the low group or high group. The experimenter stresses that this decision has to be made quickly (within 1.5 seconds). After that, subjects have to guess what urn was actually chosen. The training was split into 3 parts: estimation task, decision task and both tasks. After the instruction with slides, subjects performed 10-20 trials to familiarize themselves with the task and raised any questions they had about the task. All the training was done in one session.

Modeling

Dot stimulus

The basic model is the same as the self-consistent Bayesian observer used in Experiment 1, chapter 2 (see Appendix B). Other variants of the basic models are described below:

Logarithmic encoding: The measurement distribution and the memory recall distribution are Gaussian on the logarithm space:

$$p(m|\theta) = \frac{1}{\sigma_s \sqrt{2\pi}} e^{\frac{-(\log(m) - \log(\theta))^2}{2\sigma_s^2}} \quad (4.1)$$

$$p(m_m|m) = \frac{1}{\sigma_m \sqrt{2\pi}} e^{\frac{-(\log(m_m) - \log(m))^2}{2\sigma_m^2}} \quad (4.2)$$

Fit loss function: In this model, I use the loss function of the form $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|^n$ and fit the power n jointly with other parameters.

Fit decision prior: In this model, the observer's prior over decision variable $P('cw')$ is fit jointly with other parameters.

Symbolic number stimulus

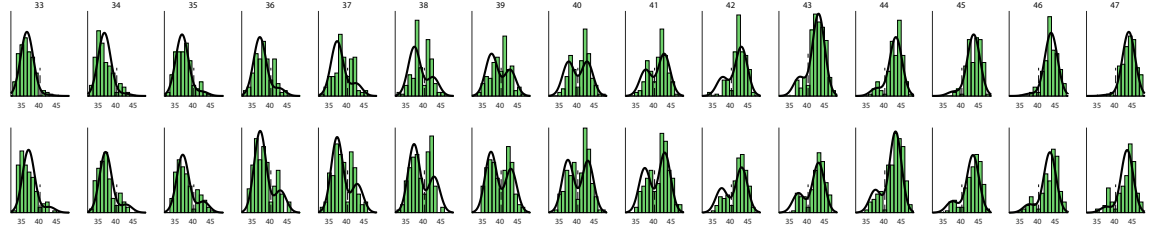
Sample Average model: The model takes the average $\hat{\mu}$ of the presented samples and reports 'greater' if $\hat{\mu}$ is greater than 40. For the estimation, it simply rounds $\hat{\mu}$ to the nearest integer.

Standard Bayesian model and self-consistent Bayesian model: The models are similar

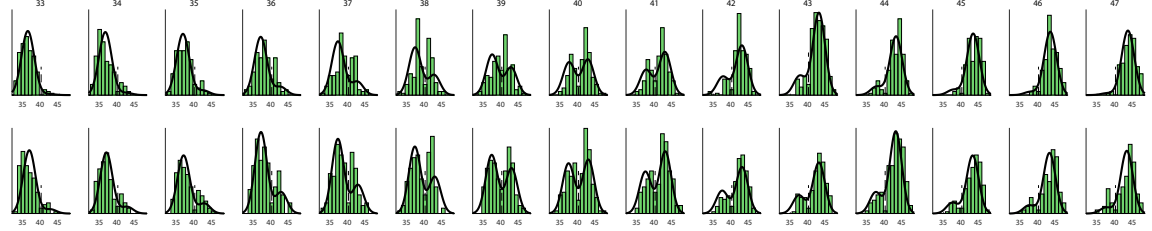
to the independent Bayesian observer and self-consistent Bayesian observer in Chapter 2 (see Methods) except that (i) all probability distribution is discrete, (ii) the prior is the experimental prior of the stimulus, (iii) the measurement distribution was the truncated discrete Gaussian described in the method above, and (iv) the self-consistent Bayesian model conditions the prior on subjects' decision in each trial.

Figure supplement

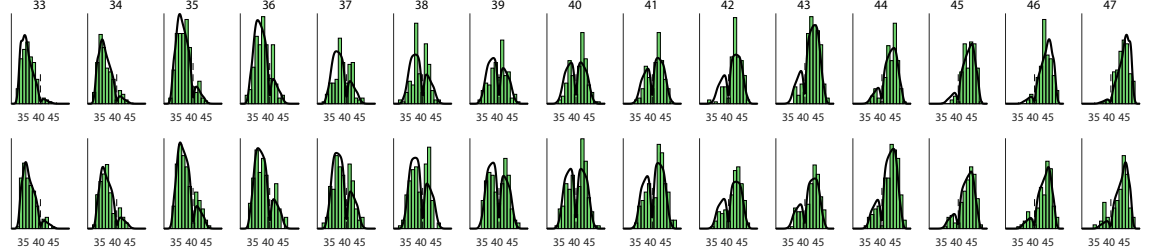
Basic model



Log encoding



Fit loss function (loss power: 0.0003)



Fit decision prior ($p(\text{"greater"}) = 0.57$)

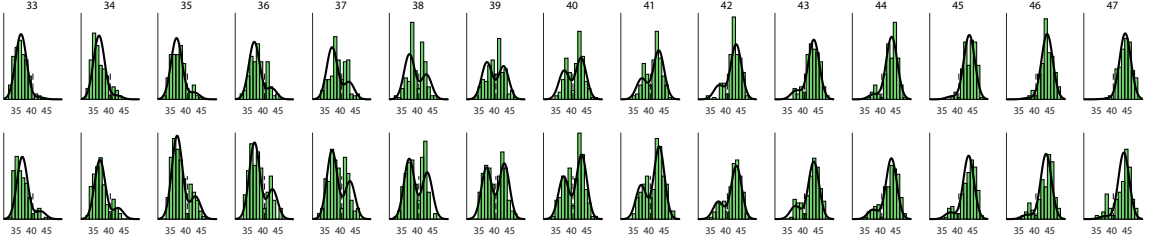


Figure 4.7: *Histogram of subjects' estimates and model fits in Experiment 1.* The green bars indicate the distribution of subjects estimates (combined subject, $n=4$) and the black curves indicate model fit. Each row corresponds to a stimulus noise level (upper: lower noise). Each column corresponds to the presented number of dots. Although the basic model can capture the general bimodal pattern in subjects' estimates, it fails to match the height and position of the distribution peak in some conditions (e.g. when the number of dots is 38). Overall, the peaks of subjects' estimate distributions are closer to the decision boundary than the model fit. The log encoding model is largely similar to the basic model with a small bias towards greater values (when the stimulus number is 40, the right peak is slightly higher than the left peak). The fit loss function model can capture the peaks in the data better than the basic and log encoding models. On the other hand, the fit decision prior model can explain the estimate bias towards greater values.

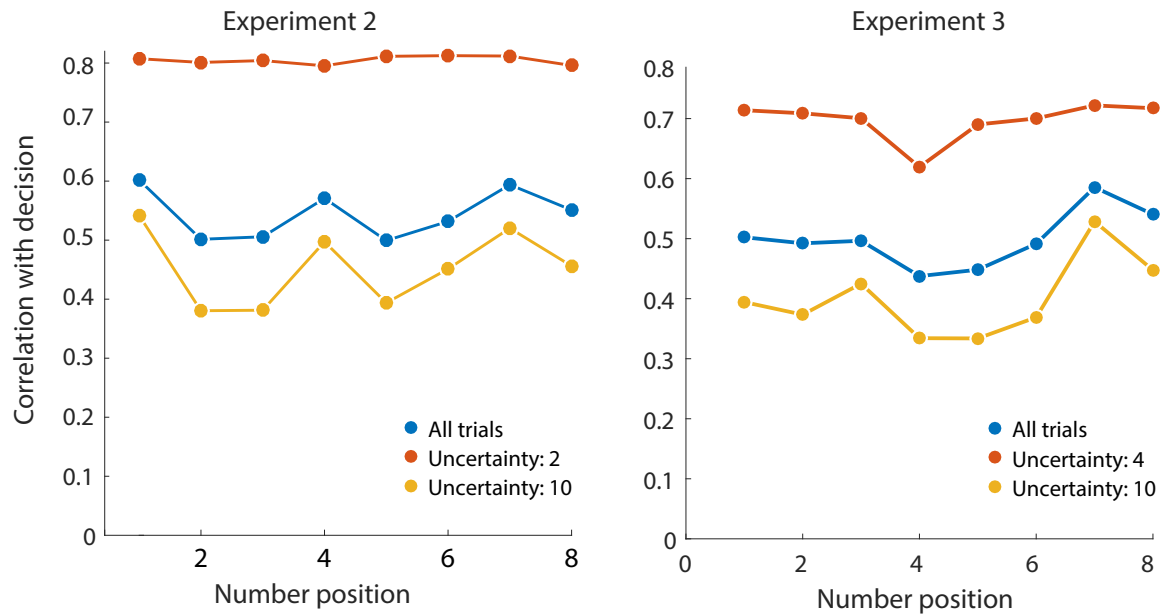


Figure 4.8: *Correlation of sample position and decision.* The graphs show Pearson correlation between subjects' decision and number position in the presented sequence. The correlation is high for low stimulus uncertainty and is almost constant across the number positions. When the stimulus uncertainty is larger, some well-known patterns emerge such as the overweight of last samples (recency effect) or first samples (primacy effect). In Experiment 2, the positional effects are not very clear. In Experiment 3, there is a significant recency effect.

Chapter 5

General discussion

5.1. Summary of thesis contribution

In chapter 2, I first replicated the bias effect in Jazayeri and Movshon (2007) with a different stimulus (lines' orientation instead of dots' motion direction, Exp. 1). Importantly, I postulated that the sequential dependency between judgments arises because humans are inclined to maintain self-consistency along the whole inference process. The self-consistency principle is formalized by self-consistent Bayesian observer model which assumes the observer treats its first decision as a fact and uses that to constrain the inference in subsequent judgments. In two additional experiments, I validated the model's key assumptions, namely, the observer performs Bayesian inference (Experiment 2) and totally trusted its first decision (Experiment 3). Combining that with model fit to data in three experiments by Zamboni et al. (2016), I demonstrated the ability of the self-consistent Bayesian observer model to quantitatively account for human behaviors across a wide range of experimental settings.

In chapter 3, I studied how subjects maintain sensory representation in working memory after the first decision. Subjects performed a sequence of decision and estimation tasks on orientation stimulus as in the previous chapter. The main difference is that subjects were given feedback on the decision task. To compare several strategies of maintaining and using sensory representation, I fit self-consistent Bayesian model to correct trials and predicting incorrect trials. The results of model comparison suggest that subjects maintain the full sensory representation and uses that to perform the estimation task.

In chapter 4, I tested the sequential dependency on a high-level stimulus variable, numerosity. In Experiment 1, a cloud of dots was briefly presented. Then subjects first indicated whether the number of dots was less or greater than 40 and subsequently re-

ported the number of dots. Experimental results show very similar bias pattern in subjects' estimates as found in the low-level stimuli (e.g. orientation). Additional experiments using symbolic number stimuli show similar characteristic biases, which suggests that sequential dependency also occurs for high-level cognitive variable and share similar features with the low-level perceptual variable.

5.2. Testable predictions of self-consistency principle

My thesis has primarily focused on the sequential tasks in which subjects perform two perceptual judgments on the same sensory evidence. Sequential dependency between judgments, however, has been found in a diverse range of other cognitive decision scenarios (see section 1.1). It is not clear whether the self-consistency principle may provide a good account of those situations as well. To demonstrate the potential of self-consistent Bayesian observer framework, I will outline in this section how we may explain sequential dependency in two other popular cognitive scenarios discussed in section 1.1. Because this may also have important implication for perceptual decision, I suggest two perceptual psychophysics experiments to test the model's prediction in analogous scenarios.

5.2.1. Sequential judgments on multiple attributes

Halo effect: In a classic study by Nisbett and Wilson (1977), subjects first made a judgment about an instructor's likability and then had to rate his specific attributes like physical attractiveness and accent. When the instructor's verbal expression was warm and respectful (highly likable), subjects' ratings of the specific attributes were strongly favorable. In contrast, the ratings were significantly downgraded when the same instructor's verbal expression was cold and distrustful. In brief, subjects' judgments on a global attribute (likability) biased judgments on other specific attributes (physical attractiveness).

Rational model: One explanation for the effect is that there exists an inherent correlational structure underlying the global attribute and the specific attributes. For instance, there is an inherent positive correlation between likability and physical attractiveness. As a result, subjects' judgments of the two attributes are naturally correlated. To demonstrate this point, I consider a toy example in which an observer makes judgments on two

attributes of an instructor: whether he is likable ($C = \{ "likable", "unlikable" \}$) and how attractive he is (Fig. 5.1a). I also assume that likability (θ) and attractiveness (γ) can be quantified on a continuous scale such that the neutral state is 0 (e.g. $\theta = -4$ indicates 'unlikable'). Crucially, the two attributes are statistically linked through the conditional probability $p(\gamma|C)$. For example, $p(\gamma|C = "likable")$ represents the distribution of physical attractiveness given that the instructor is likable. Hence if $p(\gamma|C = "likable")$ and $p(\gamma|C = "unlikable")$ are not exactly the same, there is an inherent correlation between likability and attractiveness. An example in Fig. 5.1a illustrates a plausible scenario in which attractiveness tends to be higher when the instructor is likable. In addition, we assume that the observer doesn't have access to the true value of likability and attractiveness. Due to several kinds of ambiguity (e.g. perceptual uncertainty), the observer can only obtain noisy samples m_θ and m_γ of the true likability θ and physical attractiveness γ . Given this setup, the observer's task is to use the observed samples m_θ and m_γ to make inference about likability C and attractiveness γ . I performed computer simulation for this scenario to illustrate the rational model's prediction. In the simulation, I use the probability distributions illustrated in Fig. 5.1b. Specifically, the prior distribution over C is uniform (i.e. $P(C = "likable") = P(C = "unlikable")$); the prior distributions over θ and γ are uniform in the range $[-20, 20]$. Uncertainty in the perceptual information m_θ and m_γ is characterized by $p(m_\theta|\theta)$ and $p(m_\gamma|\gamma)$ which are normal distribution centered on θ and γ , respectively. The noise of attractiveness is fixed at a moderate value (standard deviation $\sigma=7$) while likability assumes two noise levels ($\sigma_{low}=1, \sigma_{high}=18$). Note that the scale here is chosen arbitrarily; thus what matters is the relative values between parameters. The simulation results in Fig. 5.1c show qualitatively different patterns in the observer's estimate of attractiveness for different likability noise levels. When likability is less noisy, the attractiveness rating is biased away from the neutral point 0, forming a characteristic bimodal pattern. However, when likability is noisier, attractiveness rating has a more unimodal pattern centered on the true attractiveness (i.e. the diagonal line). Essentially, the rational model takes into account the noise in likability judgment when it makes judgment

about the attractiveness. On the one hand, if the likability noise is small, it will exploit the conditional dependency between likability and attractiveness to resolve the uncertainty in the attractiveness judgment. On the other hand, if the likability noise is very large, it will neglect the likability information and rely mainly on the attractiveness information. That explains the high bias, low variance in low likability noise and low bias, high variance in high likability noise. Here the bias magnitude is analogous to the correlation between attribute judgments. Nevertheless, empirical findings have shown that the correlation between attribute judgments are significantly greater than the inherent correlation of the attributes (Leuthesser et al., 1995; Thorndike, 1920). Therefore, it was hypothesized that there must be additional interaction between the judgments. This will be explored in the self-consistent model.

Self-consistent model: Built on the principle of self-consistency, I assume that after the observer makes judgment on likability, it considers the judgment as a fact and conditions the subsequent attractiveness judgment on that (Fig. 5.2a). As a result of the conditioning, the attractiveness judgments will always be biased, regardless of the likability noise and the bias pattern is analogous to the rational model’s prediction when the likability noise is zeros. The simulation result is shown in Fig. 5.2b. Importantly, if we average across all noise levels of likability, the observed biases in attractiveness judgment would be much greater than that in the rational model. That can explain the additional correlation between judgments above and beyond the intrinsic correlation between attributes.

The sequential dependency between cognitive judgments discussed above has implications for several perceptual situations. In object perception studies, there is an intrinsic hierarchical structure between global attributes (e.g. object identity) and local attributes (e.g. color, texture, shape) of an object (Kersten et al., 2004). Even in perceptual studies targeting only one low-level variable, the stimuli often contain more than one attribute. For example, research on cue combination typically involves several cues (sometimes from different sensory modalities) that subjects have to integrate to make a perceptual judgment (Ernst and Banks, 2002; Körding et al., 2007; Alais and Burr, 2004). The stimuli in

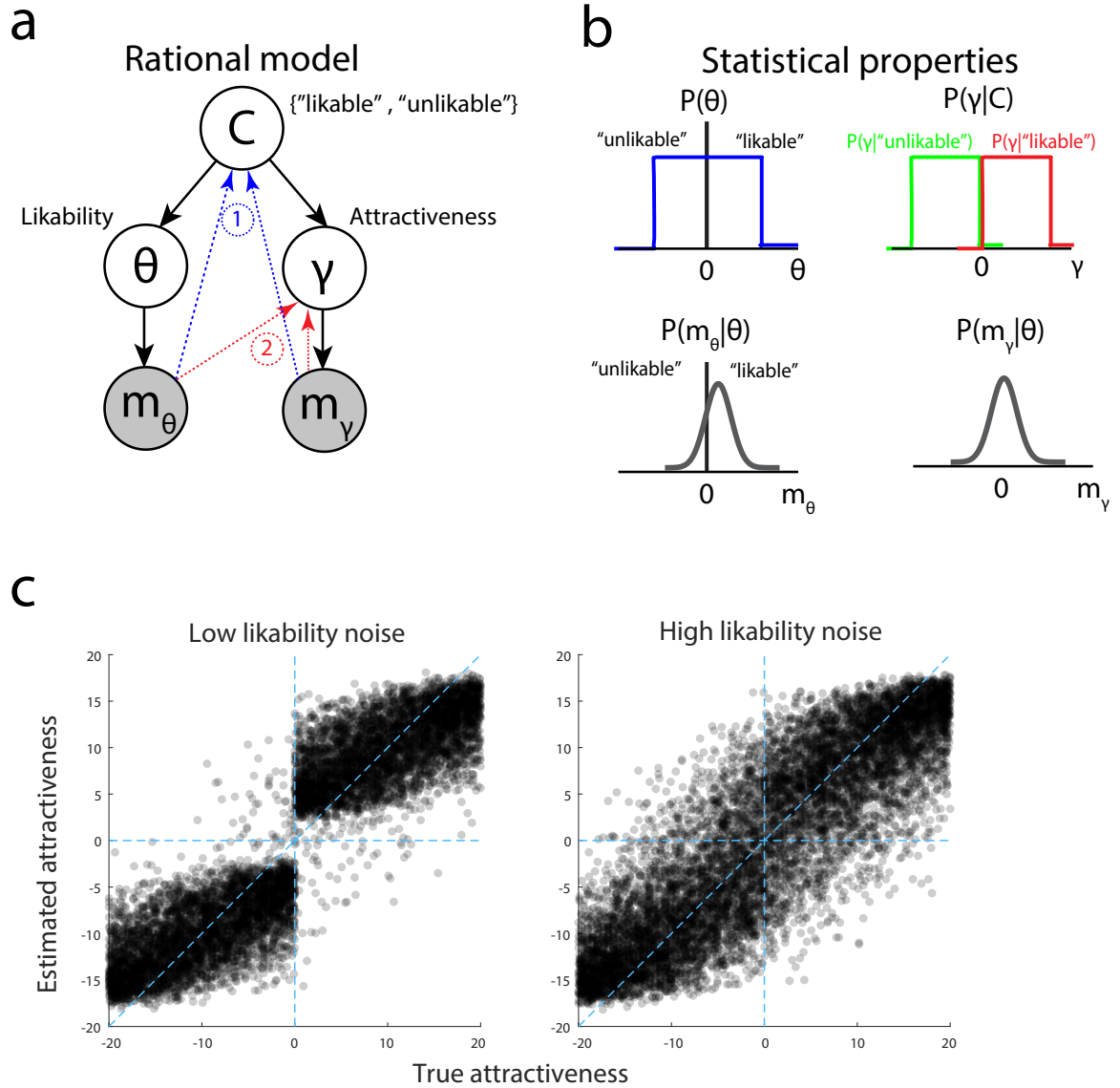


Figure 5.1: *Rational Bayesian observer of multi-attribute judgment* **(a)** Generative model: C is a categorical variable and θ is a continuous variable on likability. γ is a continuous variable on attractiveness. m_θ and m_γ are perceptual information the observer has about the instructor's likability and attractiveness. The observer uses the perceptual information to sequentially make judgments on likability and attractiveness. **(b)** There exists a correlational structure between the likability and attractiveness which is demonstrated by $P(\gamma|C)$. In this example, attractiveness is higher when the instructor is likable. **(c)** Prediction of rational model: The attractiveness judgment is plotted against the true attractiveness. When likability noise is low, attractiveness judgment is biased away from the neutral level. In contrast, the distribution of attractiveness judgment is not biased when the likability noise is high.

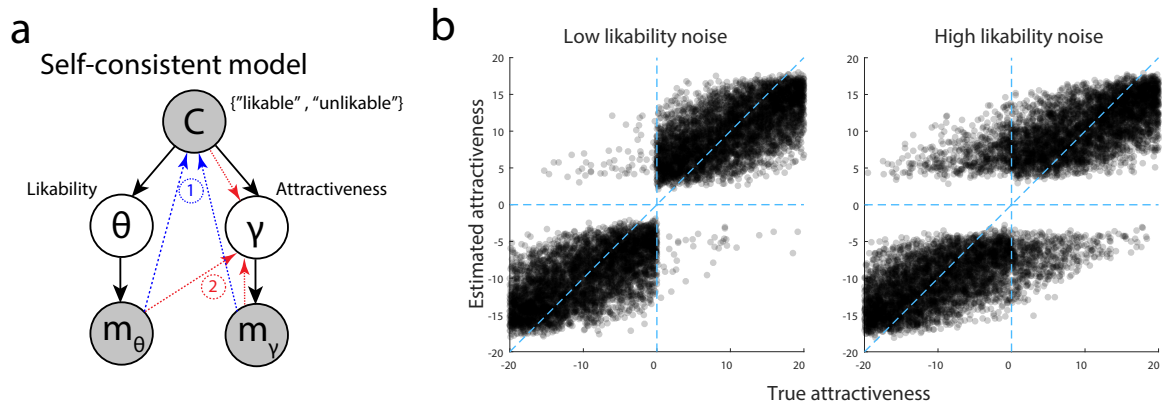


Figure 5.2: *Self-consistent observer of multi-attribute judgment* **(a)** After making judgment on likability, the observer considers the judgment \hat{C} as a fact and use that to constrain inference on attractiveness γ . **(b)** Prediction of self-consistent model: Self-consistent inference results in similar bias magnitude in attractiveness judgment across all likability noise levels. The bias magnitude is equivalent to the zero likability noise condition of the rational model. The bias patterns are different between low and high likability noise conditions because the observer is more likely to be incorrect in the likability judgment in the high noise condition which results in the more noticeable long tail in attractiveness judgment.

working memory experiments also contain multiple task-relevant features (e.g. recall the orientation of a blue line that was presented with other colored lines) (Zhang and Luck, 2008; Bays and Husain, 2008; Wilken and Ma, 2004). Moreover, as research in perception moves towards studying perceptual decision in more natural settings (e.g. more natural stimuli), it is inevitable that the perceptual situation involves stimuli with a rich hierarchical structure of multiple attributes. Therefore, it is both interesting and relevant to study whether perceptual judgments in these situations have similar sequential dependency as in cognitive judgments. In the next section, I suggest a simple perceptual experiment involving sequential decisions on two low-level attributes of the same stimulus.

Test the halo effect at low-level perception

Experimental design: I use low-level stimuli to have a better control over the attributes. The stimuli and experimental set-up are illustrated in Fig. 5.3. The stimulus is an array of colored line segments. I take the overall orientation of the array as the global attribute θ . Two noise levels of stimulus orientation can be obtained by changing the variance of the distribution the lines' orientation are sampled from. The hue of line segments represents the

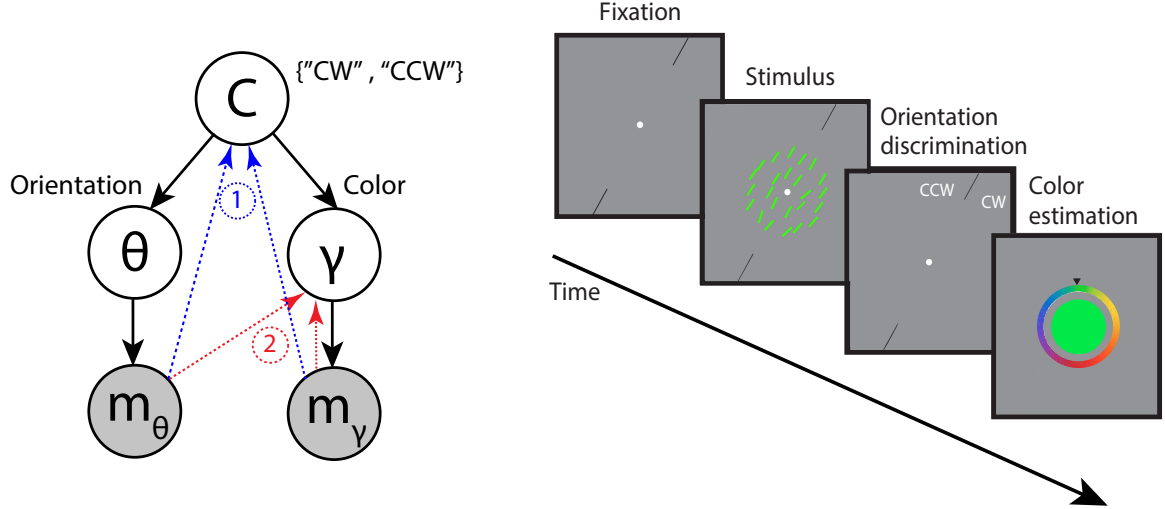


Figure 5.3: *Experimental test of halo effect at low-level perception.* Subjects are presented with a color array of line segments. They first make a categorical decision on the overall orientation of the array. Subsequently, they rotate an arrow around the color wheel to match the perceived color of the array.

specific attribute γ . To obtain a moderate noise level of the hue attribute, we may either tune the chromaticity of the background or the saturation of the line segments. The joint distribution of the orientation and hue dimensions is the same as in the above simulation. After viewing the stimulus, subjects first indicate whether the overall orientation of the array was clockwise or counter-clockwise of a reference orientation. Then they have to estimate the hue by adjusting an arrow around a color wheel. The color that the arrow points to is presented on a disk inside the color wheel.

Analysis: When the hue values are close to the center range, the rational model predicts larger repulsive biases in color estimate for lower orientation noise level. In contrast, the self-consistent model predicts the same biases in color estimate across all orientation noise levels. As a simple statistical test, we can conduct a two-tailed t-test between the low and high orientation noise conditions with the color judgment bias as the dependent variable. Alternatively, we can fit the two models to subjects' data to see which model can result in a better fit. Because the two models have the same generative model and parameters, I can directly compare the maximum log likelihood of model fit.

Anticipated and alternative result: If the results favor the self-consistent observer model,

it demonstrates that the halo effect appears even at low-level perception. Alternatively, we may find lower biases for higher orientation noise as predicted by the rational model. In this scenario, it could be that at low-level perception, humans are able to take into account the dependencies between attributes to make a rational judgment. However, it does not rule out the possibility that in a more complex setting, people will eventually revert to self-consistent inference. Another explanation for the null result is that subjects fail to learn the experimentally induced probabilistic relations between features. Therefore, it is essential to have sufficient training and manipulation check to guarantee that subjects learn the stimulus statistics.

5.2.2. Incorporating additional evidence after a preliminary decision

Sunk cost fallacy: Another kind of bias that may arise in sequential decisions is when we first make a preliminary decision on the current evidence and then make another decision given additional evidence. An example experiment in real-life setting was described in Arkes and Blumer (1985). At a theater, the ticket price was randomly chosen between a normal (\$15) and a discount price (\$2) for the buyers within a certain period of time. The tickets were marked so that the experimenter could track how many people would go to the show for each ticket price. The result shows that more people with the normal-priced tickets went to the show. If we consider the act of buying ticket as the first decision and the act of going to the show as the second decision, the findings suggest that subjects' second decision was biased by how much they had decided to spend on the ticket (i.e. the sunk cost).

In a similar decision situation but with perceptual stimuli, Bronfman et al. (2015) found that after having made a preliminary decision, subjects tended to underweight the subsequent evidence in the final judgment. That results in the bias effect like the sunk-cost fallacy. Along the same line, studies on perceptual confidence suggest that after making a decision, subjects still accumulate sensory evidence until a confidence judgment is made (Pleskac and Busemeyer, 2010) and there exist biases in the post-decision evidence accumulation (Yu et al., 2015a; Navajas et al., 2016). Another relevant line of works in perception posits that there are two pathways in visual processing: a fast, low spatial frequency pathway

(magnocellular) and a slow, high spatial frequency pathway (parvocellular) (Bar, 2004; Bar et al., 2006; O’Callaghan et al., 2016). The authors suggest that the fast pathway comes up with a preliminary coarse judgment (e.g. on the context) and uses that to constrain and facilitate the final judgment when the fine-detailed sensory information comes in from the slow pathway. Those perceptual situations bear a close resemblance to the sunk cost fallacy scenario in cognitive judgments. It is possible that the biases in perceptual studies share similar mechanisms as the biases in cognitive literature. Therefore, I propose a perceptual experiment with similar task setting to test whether the self-consistency principle may account for those biases.

Experimental design: The experimental setup is similar to Experiment 1 in chapter 2. In each trial, I draw a stimulus orientation from a uniform distribution around a randomly chosen reference orientation. Then I draw two samples of line segments, each sample containing 24 segments, from a Gaussian centered on the stimulus orientation. Subjects are first presented with one sample and have to tell whether the stimulus orientation is CCW or CW of the reference. After that, they view the second sample and have to estimate the stimulus orientation by adjusting a probe line. In this experiment, it is important to instruct subjects about the structure of the stimulus and the task. Specifically, subjects are instructed that in each trial, there is a box containing a large number of line segments. First, a set of line segments are drawn from this box and they have to guess whether the mean orientation of all line segments in the box was CCW or CW of the reference. Then I show them another set drawn from the same box and they have to estimate the mean orientation of all line segments in the box.

Prediction of alternative models: In the simulation, I assume the prior distribution of stimulus orientation is uniform in the range $[-21, 21]$ deg; sensory noise follows a normal distribution which is centered on the stimulus orientation and has a standard deviation of 7 deg. Fig. 5.4a shows the generative model of the task and the inference process of the rational model. The observer uses the first sensory sample m_1 to decide whether the stimulus orientation is clockwise or counterclockwise relative to the reference. Then upon

receiving the second sample m_2 they use both to estimate the stimulus orientation. The rational model predicts no repulsive biases in the estimation; thus the estimation pattern is unimodally distributed around the true stimulus orientation (Fig. 5.4a). In contrast, the self-consistent model treats the first decision as a fact and uses that in the subsequent estimation with both samples. That leads to repulsive biases and accordingly bimodal pattern of estimate distribution (Fig. 5.4b). The most notable trials is when the first sample leads to the wrong conclusion and the second sample provides contradictory information (Fig. 5.4c). The rational observer can use the second sample to revise the judgment and make unbiased estimates. On the contrary, the self-consistent observer distorts the sensory evidence to make it consistent with the first decision.

5.3. Self-consistent inference

5.3.1. Benefits of self-consistency

So far the self-consistent Bayesian observer model has been shown to account for several post-decision biases across different stimuli and situations. However, it is unclear what benefit the self-consistent inference may bring about in those situations. In the absence of memory noise, it is quite straightforward that a standard Bayesian observer should outperform the self-consistent observer in terms of accuracy. With memory noise, I have shown in chapter 2 that subjects' sequential judgments are consistent with each other only when they perform self-consistent inference (see section 2.5). Nevertheless, it still doesn't directly address any advantage the self-consistent inference may have. To target this issue, I took the fit parameters of nine subjects in Exp. 1-3, chapter 2 and performed model simulation to compute the mean squared error for the self-consistent Bayesian model and the standard Bayesian model. Interestingly, the self-consistent model outperforms the standard Bayesian model for 6 out of 9 subjects (Fig. 5.5a). The results are similar regardless of whether I take into account the motor noise in the computation or not. A closer look at the mean squared error indicates a pretty consistent pattern across subjects that the lower the sensory noise is, the better the self-consistent model is compared to the standard Bayesian model. In fact, this pattern is not solely modulated by the magnitude of sensory noise itself but the

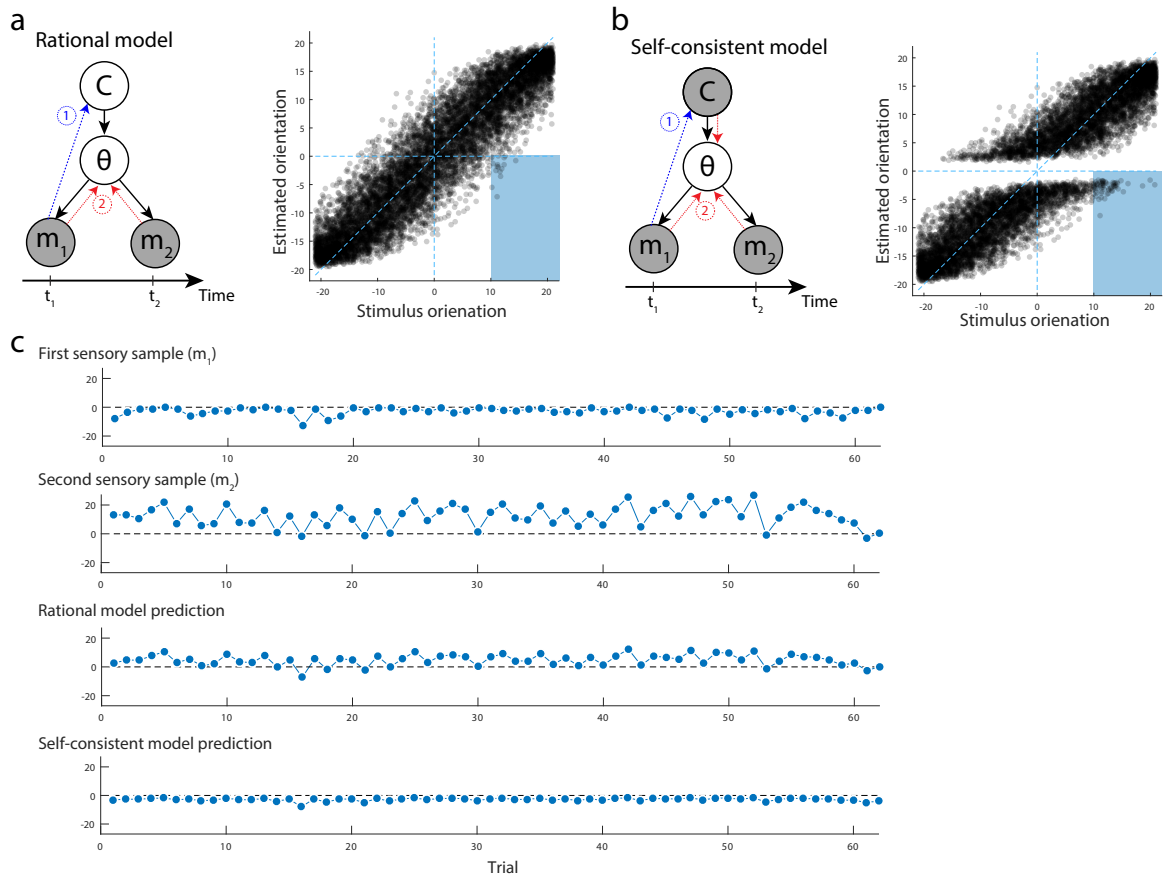


Figure 5.4: *Modeling sunk cost fallacy* (a) Rational observer model: The observer uses the first sensory evidence m_1 to make a categorical decision $C = \{“CW”, “CCW”\}$. After receiving the second sensory evidence m_2 , it uses both pieces of evidence to make inference on the stimulus orientation θ . (b) Self-consistent observer model: The first categorical decision is the same as the rational model. Importantly, the observer considers that decision as completely true and uses that decision to constrain inference in the following estimation. As a result, the distribution of estimates has a characteristic bimodal pattern. (c) Distortion of sensory evidence: To illustrate how the self-consistent observer distorts the evidence, I picked the trials for which the stimulus orientation is in the range $[10, 21]$ and the first sensory evidence falls on the incorrect side, that is $m_1 < 0$ (the blue-shaded region in b). Because the stimulus orientation is far away from the boundary, it is very likely that the second sensory sample falls on the correct side. Therefore a rational model that merely integrates the two samples will be likely to make estimates that fall on the correct side. In contrast, the self-consistent observer maintains self-consistent with the first decision by distorting the sensory evidence. That results in the estimates that always fall on the incorrect side like the first sensory evidence.

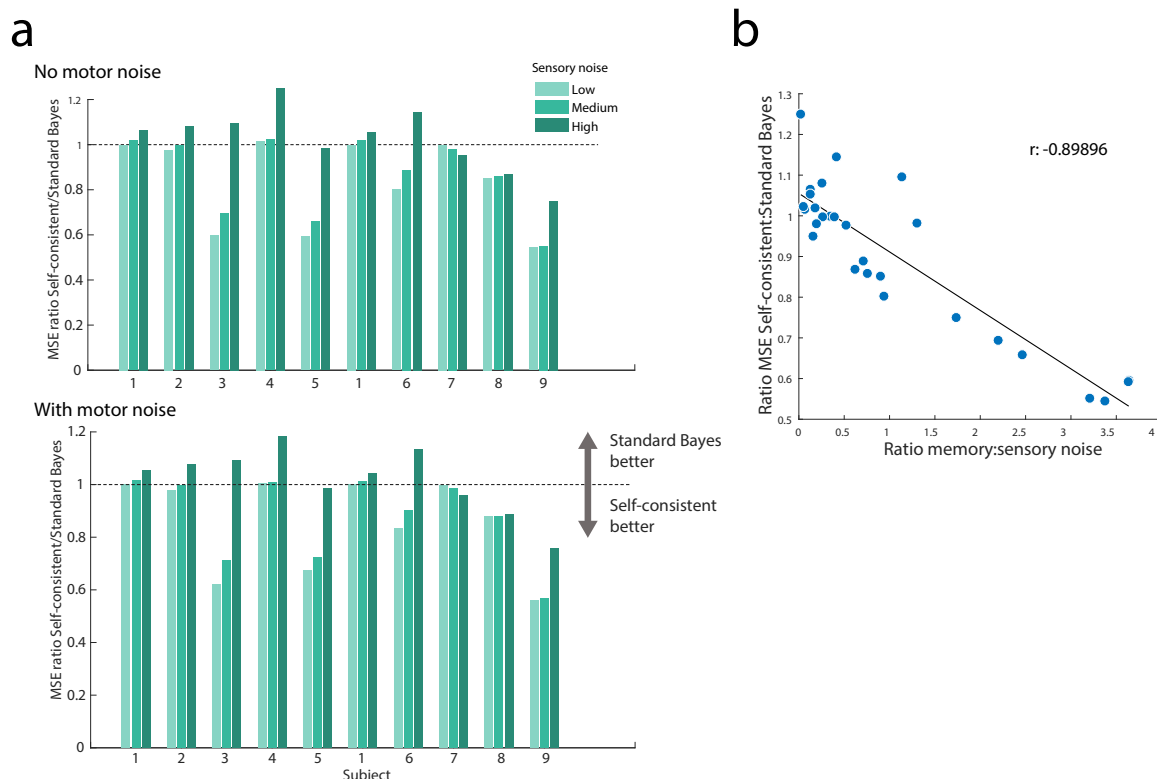


Figure 5.5: *Advantage of self-consistent Bayesian observer* (a) Ratio of mean squared error (MSE) between self-consistent and standard Bayesian model are shown for each individual subject and stimulus noise level. A ratio less than 1 means the self-consistent Bayesian model is better. The mean squared error is computed across the full range of stimulus orientation. The upper graph shows the results when motor noise is not taken into account and in the lower graph, the motor noise is incorporated in the error computation. (b) Correlation between the ratio of MSE and the ratio of memory to sensory noise. Each dot corresponds to one subject and one noise level.

sensory noise relative to memory noise. There is a high negative correlation between the model MSE ratio and the sensory to memory noise ratio (Fig. 5.5b). More specifically, the higher the memory noise is relative to sensory noise, the better the self-consistent model is compared to the standard Bayesian model. This is consistent with simulation results documented elsewhere (Luu et al., 2017).

To understand why self-consistent inference is beneficial when memory noise is high compared to sensory noise, let's consider an extreme situation when the sensory noise is negligible and memory noise is almost infinite. In that case, the first decision based on

initial sensory information is very accurate whereas the memory recall contains insignificant information about the stimulus. As a result, the standard Bayesian model that only uses memory recall to make estimation will be worse than the self-consistent Bayesian model that uses both the preceding decision and the memory recall. Here we assume that the categorical decision is not degraded by memory noise whereas the sensory representation is. The assumption is plausible because the one-bit representation of binary decision is more robust to noise than the continuous sensory representation of the stimulus. Note that the above analysis is only limited to the simple experiments in Chapter 2. Future works are needed to determine whether the findings hold for other more natural scenarios and more interestingly, whether humans can adaptively adjust the inference strategy to maximize performance across different situations.

Another potential benefit of self-consistent inference is that it can save computational costs in situations where the task sequence is more complex than the simple experiments in laboratory settings. Starting from the top, it reduces the decision tree at every level of the hierarchy by considering only the chosen branch, substantially reducing the overall computational complexity and cost associated with solving the inference problem. Self-consistent inference may represent a general strategy for the brain to address the cost-accuracy trade-off when solving hierarchical decision-making problems. This also may have important implications for learning and belief updating in biological as well as artificial neural networks, in particular for networks that are aimed at learning a generative model (*e.g.* deep belief networks).

5.3.2. Conditioning likelihood and loss function

In my thesis, the self-consistent Bayesian observer enforces self-consistency by conditioning the prior on the preceding decision. Although the model can account well for the data, it is not the only way to implement self-consistency in the Bayesian framework. In principle, we can also condition the likelihood function and/or the loss function to make it consistent with the preceding decision. The interpretation, however, would be different for conditioning different components of a Bayesian observer. Conditioning the prior means changing the

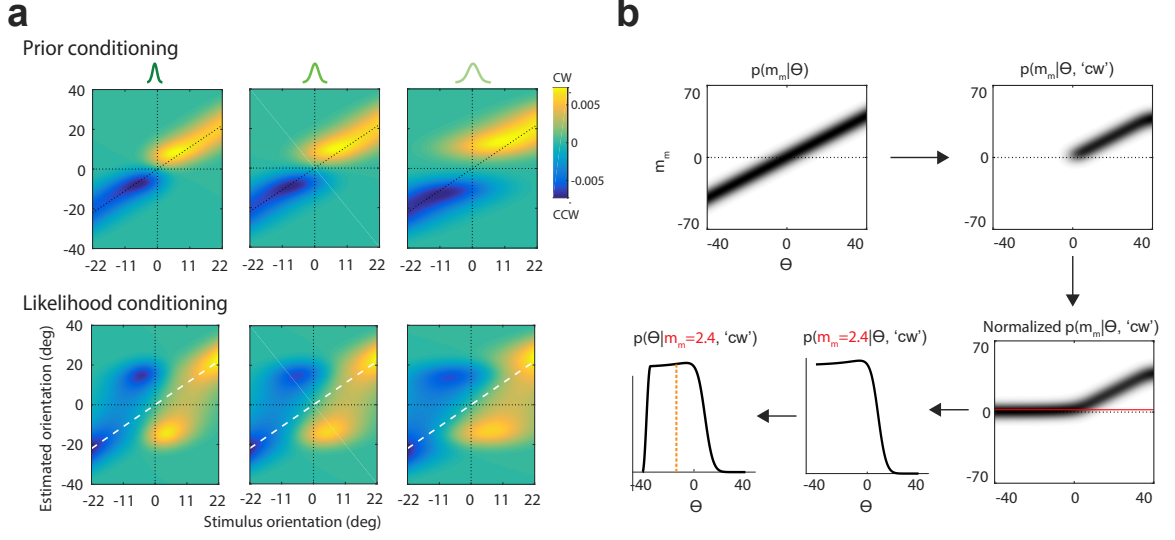


Figure 5.6: *Conditioning the likelihood.* (a) The distribution of the observer's estimate is shown when conditioning occurs in the prior (upper panel) or the likelihood (lower panel). Three columns correspond to three stimulus noise levels with increasing noise from left to right. The parameters used for the simulation are taken from the fit parameters of combined subject in Exp. 1, chapter 2. (b) Likelihood conditioning flips the estimate. $p(m_m|\theta)$ indicates the distribution of memory recall m_m given a stimulus orientation θ . The conditioned distribution $p(m_m|\theta, 'cw')$ indicates the distribution of memory recall that results from the sensory measurement m on the *cw* side. After normalization so that $\int p(m_m|\theta, 'cw') dm = 1$, the likelihood function is skewed towards the *ccw* side. As a result, the estimate from the posterior $p(\theta|m_m, 'cw')$ is on the *ccw* part as illustrated in the case $m_m = 2.4$. Note that this just happens for the samples m_m that are around 0 which explains why the flipping only occurs for the stimulus orientation close to the boundary.

observers' belief about the statistical regularity of the external world. On the other hand, conditioning the likelihood means changing the observer's belief about the internal sensory processing. By conditioning the loss function, the observer changes belief about how each action is rewarded.

In terms of the observed behavior, different types of conditioning may have very distinct signatures. To illustrate the difference between prior and likelihood conditioning, I computed the prediction of the self-consistent observer for the basic task in Exp. 1, chapter 2 in which subjects made a decision and then an estimation of the stimulus orientation (Fig. 5.6a). The patterns of estimate distribution are markedly different between these two types of conditioning. Likelihood conditioning leads to counterintuitive results in which the

observer’s estimates are contradictory with the decision when the stimulus orientation is close to the decision boundary (see Fig. 5.6b for technical detail). It also clearly shows that for the experimental settings studied in my thesis, prior conditioning provides a better explanation than the likelihood conditioning. That makes sense because conceptually, changing the prior and the loss function seems to be the more intuitive and plausible than the likelihood because the sensory process is often more stable and more accessible to the observer than the external world and the reward structure. As a result, likelihood conditioning may happen only in special cases such as patients with neurological disorders that makes the sensory processing unstable. Computationally, conditioning the loss function would be more complicated because the observer has to take into account the motor noise (Tank and Stocker, 2014) and the reward contingency of each action (Trommershäuser et al., 2003, 2008). Note that the different types of conditioning are not mutually exclusive. In other words, the self-consistent observer can apply a combination of those conditioning strategies. For example, both the prior and the likelihood can be conditioned simultaneously which indicates that the observer attempts to be self-consistent in both encoding and decoding.

5.4. General implications for decision-making study

The fact that subjects condition their estimate on their preceding decision does not imply that they are necessarily fully confident in their decision; we propose that they simply do so in order to remain self-consistent. Our results show that conditioning is statistically independent of the difficulty and thus on subjects’ confidence in their discrimination judgment (*i.e.* their psychometric function). However, it remains an interesting open question particularly in context of the ongoing discussion about decision confidence (Kepecs et al., 2008; van den Berg et al., 2016; Fleming and Daw, 2017) whether or not conditioning improves subjects’ confidence in their subsequent estimate since it leads to a reduced posterior distribution.

Another interesting question is whether an explicit categorical commitment is necessary to induce self-consistent conditioning or whether the brain always, and thus implicitly, performs conditioned inference (Ding et al., 2017). This question is difficult to answer because

without explicit access to a subject’s trial-by-trial categorical judgment, differences in the subsequent feature inference process are often hard to distinguish statistically. Only in special cases when, for example, a clear discrimination boundary is present, these differences have a clear behavioral signature that can be identified (*e.g.* the characteristic repulsive estimation pattern in Experiment 1, see Figure 2.1). Identification is further complicated because repulsive biases may also arise due to other effects such as the efficient adaptation of sensory encoding resources (Wei and Stocker, 2015), which likely takes place during perceptual learning (*e.g.* Szpiro et al. (2014); see also Wei and Stocker (2017)). In fact, it may well be that implicit self-consistent inference is the fundamental process by which the brain solves inference problems, yet its behavioral characteristics are simply not often apparent. For example, we expect self-consistent conditioning to happen implicitly in object recognition of learned object categories: when an observer recognizes an object as *e.g.* an “apple”, the percept of the object’s features (*e.g.* the color) will automatically be conditioned on that recognized category. In order to detect the effects of this conditioning in perceptual behavior, however, we would need to know the specifics of the learned generative models over the object categories, which is typically a difficult problem. In other situations, such as in a typical psychophysical experiment with its sparse and artificial stimuli and little context, the observer may simply be given little incentive to interpret the stimulus within a hierarchical representation (generative model). Because the self-consistent inference model over a flat generative model is identical to an optimal Bayesian observer model, the large number of studies that have shown that perception is well explained as optimal Bayesian inference may actually not be conclusive; their data is equally well explained by the self-consistent inference model! This is obviously a strong hypothesis that needs further experimental evaluation.

Results from recent physiological recordings in primates suggest not only that decision-making is associated with rapid cortical state-changes (Meindertsma et al., 2017) but also that decision-related signals are fed back along the perceptual processing pathway all the way to early sensory areas (Nienborg and Cumming, 2009; Siegel et al., 2015). The proposed

self-consistent observer model provides a novel computational interpretation of these neural signals: at the moment a decision is reached the belief state of the brain rapidly changes (in favor of the choice made, expressing strong over-confidence, see Peters et al. (2017)) and is fed back to ensure that the perceptual inference process remains consistent across the different cortical levels of representation at any moment in time. The self-consistent model may prove a useful hypothesis to constructively explore the function and purpose of such decision-related signal flows in the brain. Future work needs to explore how exactly our model formulation can be interpreted at a more mechanistic neural level (Luu and Stocker, 2016). Various theoretical frameworks have been proposed for how the brain might perform Bayesian inference (*e.g.* Ma et al. (2006); Simoncelli (2009); Wei and Stocker (2012); Pitkow and Angelaki (2017)). It remains an interesting challenge to investigate how these frameworks can incorporate the self-consistency constraint that we propose here, in particular the process of quickly and flexibly imposing a conditional prior.

Our results show interesting parallels to many well-known bias phenomena in cognition and economics, such as confirmation bias (Nickerson, 1998), biases associated with cognitive consistency (Brehm, 1956; Abelson, 1968) and dissonance (Festinger, 1957; Festinger and Carlsmith, 1959; Egan et al., 2010; Sharot et al., 2010), as well as loss aversion and the sunk cost fallacy (Kahneman and Tversky, 1984). Our findings seem also aligned with results from human probability judgments over hierarchical representations which found that subjects rather follow individual probability branches than to resolve the entire probability tree (Lagnado and Shanks, 2003). It will be interesting to explore to what degree the proposed self-consistent model generalizes to these cognitive phenomena and is able to provide a parsimonious, quantitative explanation.

Bibliography

- Abelson, R. P., editor (1968). *Theories of cognitive consistency; a sourcebook*. Rand McNally, Chicago.
- Abrahamyan, A., Silva, L. L., Dakin, S. C., Carandini, M., and Gardner, J. L. (2016). Adaptable history biases in human perceptual decisions. *Proceedings of the National Academy of Sciences*, 113(25):E3548–E3557.
- Acerbi, L., Wolpert, D. M., and Vijayakumar, S. (2012). Internal representations of temporal statistics and feedback calibrate motor-sensory interval timing. *PLoS computational biology*, 8(11):e1002771.
- Akaishi, R., Umeda, K., Nagase, A., and Sakai, K. (2014). Autonomous mechanism of internal choice estimate underlies decision inertia. *Neuron*, 81(1):195–206.
- Alais, D. and Burr, D. (2004). The ventriloquist effect results from near-optimal bimodal integration. *Current biology*, 14(3):257–262.
- Angela, J. Y. and Cohen, J. D. (2009). Sequential effects: superstition or rational behavior? In *Advances in neural information processing systems*, pages 1873–1880.
- Anobile, G., Cicchini, G. M., and Burr, D. C. (2016). Number as a primary perceptual attribute: A review. *Perception*, 45(1-2):5–31.
- Anobile, G., Turi, M., Cicchini, G. M., and Burr, D. C. (2015). Mechanisms for perception of numerosity or texture-density are governed by crowding-like effects. *Journal of Vision*, 15(5):4–4.
- Arkes, H. R. and Blumer, C. (1985). The psychology of sunk cost. *Organizational behavior and human decision processes*, 35(1):124–140.
- Arons, L. and Irwin, F. W. (1932). Equal weights and psychophysical judgments. *Journal of Experimental Psychology*, 15(6):733.
- Arrighi, R., Togoli, I., and Burr, D. C. (2014). A generalized sense of number. *Proceedings of the Royal Society of London B: Biological Sciences*, 281(1797):20141791.
- Ascher, D. and Grzywacz, N. M. (2000). A bayesian model for the measurement of visual velocity. *Vision research*, 40(24):3427–3434.
- Ayoub, K. and Payne, K. (2016). Strategy in the age of artificial intelligence. *Journal of Strategic Studies*, 39(5-6):793–819.
- Bar, M. (2004). Visual objects in context. *Nature Reviews Neuroscience*, 5(8):617.
- Bar, M., Kassam, K. S., Ghuman, A. S., Boshyan, J., Schmid, A. M., Dale, A. M., Hämäläinen, M. S., Marinkovic, K., Schacter, D. L., Rosen, B. R., et al. (2006). Top-down facilitation of visual recognition. *Proceedings of the national academy of sciences*, 103(2):449–454.

- Bays, P. M., Gorgoraptis, N., Wee, N., Marshall, L., and Husain, M. (2011). Temporal dynamics of encoding, storage, and reallocation of visual working memory. *Journal of Vision*, 11(10):6.
- Bays, P. M. and Husain, M. (2008). Dynamic shifts of limited working memory resources in human vision. *Science*, 321(5890):851–854.
- Bem, D. J. (1972). Self-perception theory¹. In *Advances in experimental social psychology*, volume 6, pages 1–62. Elsevier.
- Blakemore, C. (1993). *Vision: Coding and efficiency*. Cambridge University Press.
- Brainard, D. H. and Freeman, W. T. (1997). Bayesian color constancy. *JOSA A*, 14(7):1393–1411.
- Brainard, D. H., Longere, P., Delahunt, P. B., Freeman, W. T., Kraft, J. M., and Xiao, B. (2006). Bayesian model of human color constancy. *Journal of Vision*, 6(11):10–10.
- Brainard, D. H. and Vision, S. (1997). The psychophysics toolbox. *Spatial vision*, 10:433–436.
- Braun, A., Urai, A. E., and Donner, T. H. (2018). Adaptive history biases result from confidence-weighted accumulation of past choices. *Journal of Neuroscience*, pages 2189–17.
- Brehm, J. (1956). Postdecision changes in the desirability of alternatives. *Journal of Abnormal and Social Psychology*, 52(3):384ff.
- Brezis, N., Bronfman, Z. Z., and Usher, M. (2015). Adaptive spontaneous transitions between two mechanisms of numerical averaging. *Scientific reports*, 5:10415.
- Bronfman, Z. Z., Brezis, N., Moran, R., Tsetsos, K., Donner, T., and Usher, M. (2015). Decisions reduce sensitivity to subsequent information. *Proceedings of the Royal Society of London B: Biological Sciences*, 282(1810).
- Brunswik, E. and Kamiya, J. (1953). Ecological cue-validity of ‘proximity’ and of other gestalt factors. *The American journal of psychology*, 66(1):20–32.
- Burge, J., Fowlkes, C. C., and Banks, M. S. (2010). Natural-scene statistics predict how the figure-ground cue of convexity affects human depth perception. *Journal of Neuroscience*, 30(21):7269–7280.
- Burr, D. and Ross, J. (2008). A visual sense of number. *Current biology*, 18(6):425–428.
- Burr, D. C., Ross, J., Binda, P., and Morrone, M. C. (2010). Saccades compress space, time and number. *Trends in cognitive sciences*, 14(12):528–533.
- Busse, L., Ayaz, A., Dhruv, N. T., Katzner, S., Saleem, A. B., Schölvink, M. L., Zaharia, A. D., and Carandini, M. (2011). The detection of visual contrast in the behaving mouse. *Journal of Neuroscience*, 31(31):11351–11361.

- Cheadle, S., Wyart, V., Tsetsos, K., Myers, N., De Gardelle, V., Castañón, S. H., and Summerfield, C. (2014). Adaptive gain control during human perceptual choice. *Neuron*, 81(6):1429–1441.
- Chen, M. and Risen, J. (2010). How choice affects and reflects preferences: Revisiting the free-choice paradigm. *Journal of Personality and Social Psychology*, 99(4):573–594.
- Cicchini, G. M., Anobile, G., and Burr, D. C. (2016). Spontaneous perception of numerosity in humans. *Nature Communications*, 7:12536.
- Cohen, J. B. and Goldberg, M. E. (1970). The dissonance model in post-decision product evaluation. *Journal of Marketing Research*, pages 315–321.
- Collier, G. (1951). An investigation of the independence of successive responses for behavior at the visual threshold. *Amer. Psychol*, 6:278.
- Collier, G. (1954). Intertrial association at the visual threshold as a function of intertrial interval. *Journal of Experimental Psychology*, 48(5):330.
- Collier, G. and Verplanck, W. S. (1958). Nonindependence of successive responses at the visual threshold as a function of interpolated stimuli. *Journal of experimental psychology*, 55(5):429.
- Cornsweet, T. (2012). *Visual perception*. Academic press.
- Davelaar, E. J., Goshen-Gottstein, Y., Ashkenazi, A., Haarmann, H. J., and Usher, M. (2005). The demise of short-term memory revisited: empirical and computational investigations of recency effects. *Psychological review*, 112(1):3.
- Day, W. F. (1956). Serial non-randomness in auditory differential-thresholds as a function of interstimulus interval. *The American journal of psychology*, 69(3):387–394.
- Day, W. F. (1957). A descriptive analysis of the non-randomness of serial threshold-data. *The American journal of psychology*, 70(2):238–242.
- de Lange, F. P., Rahnev, D. A., Donner, T. H., and Lau, H. (2013). Prestimulus oscillatory activity over motor cortex reflects perceptual expectations. *Journal of Neuroscience*, 33(4):1400–1410.
- Ding, S., Cueva, C., Tsodyks, M., and Qian, N. (2017). Visual perception as retrospective bayesian decoding from high- to low-level features. *Proc. National Academies of Sciences U.S.A.*, Early edition.
- Dorfman, D. D. and Biderman, M. (1971). A learning model for a continuum of sensory states. *Journal of Mathematical Psychology*, 8(2):264–284.
- Drugowitsch, J., Wyart, V., Devauchelle, A.-D., and Koechlin, E. (2016). Computational precision of mental inference as critical source of human choice suboptimality. *Neuron*, 92(6):1398–1411.

- Egan, L., Bloom, P., and Santos, L. (2010). Choice-induced preferences in the absence of choice: Evidence from a blind two choice paradigm with young children and capuchin monkeys. *Journal of Experimental Social Psychology*, 46:2014–207.
- Ernst, M. O. and Banks, M. S. (2002). Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415(6870):429–433.
- Ernst, M. O., Banks, M. S., and Bühlhoff, H. H. (2000). Touch can change visual slant perception. *Nature neuroscience*, 3(1):69.
- Farnsworth, D. (1943). The farnsworth-munsell 100-hue and dichotomous tests for color vision. *JOSA*, 33(10):568–578.
- Fechner, G. (1860). *Elemente der psychophysik*. Leipzig: Breitkopf & Hartel.
- Fernberger, S. W. (1920). Interdependence of judgments within the series for the method of constant stimuli. *Journal of Experimental Psychology*, 3(2):126.
- Festinger, L. (1957). *Theory of Cognitive Dissonance*. Stanford University Press, Stanford, CA.
- Festinger, L. and Carlsmith, J. (1959). Cognitive consequences of forced compliance. *Journal of Abnormal and Social Psychology*, 58(2):203ff.
- Festinger, L., Riecken, H. W., and Schachter, S. (1956). *When prophecy fails*. Minneapolis, MN, US: University of Minnesota Press.
- Fleming, S. and Daw, N. (2017). Self-evaluation of decision-making: A general Bayesian framework for metacognitive computation. *Psychological Review*, 124(1):91–114.
- Frenkel, O. J. and Doob, A. N. (1976). Post-decision dissonance at the polling booth. *Canadian Journal of Behavioural Science/Revue canadienne des sciences du comportement*, 8(4):347.
- Fründ, I., Wichmann, F. A., and Macke, J. H. (2014). Quantifying the effect of intertrial dependence on perceptual decisions. *Journal of vision*, 14(7):1–9.
- Gebuis, T. and Reynvoet, B. (2011). Generating nonsymbolic number stimuli. *Behavior research methods*, 43(4):981–986.
- Gebuis, T. and Reynvoet, B. (2012). The interplay between nonsymbolic number and its continuous visual properties. *Journal of Experimental Psychology: General*, 141(4):642.
- Geisler, W. S. and Perry, J. S. (2009). Contour statistics in natural images: Grouping across occlusions. *Visual neuroscience*, 26(1):109–121.
- Geisler, W. S., Perry, J. S., Super, B., and Gallogly, D. (2001). Edge co-occurrence in natural images predicts contour grouping performance. *Vision research*, 41(6):711–724.
- Gifford, A. M., Cohen, Y. E., and Stocker, A. A. (2014). Characterizing the impact of category uncertainty on human auditory categorization behavior. *PLoS computational biology*, 10(7):e1003715.

- Girshick, A. R., Landy, M. S., and Simoncelli, E. P. (2011). Cardinal rules: visual orientation perception reflects knowledge of environmental statistics. *Nature neuroscience*, 14(7):926.
- Glaze, C. M., Kable, J. W., and Gold, J. I. (2015). Normative evidence accumulation in unpredictable environments. *Elife*, 4:e08825.
- Gold, J. I., Law, C.-T., Connolly, P., and Bennur, S. (2008). The relative influences of priors and sensory evidence on an oculomotor decision variable during perceptual learning. *Journal of Neurophysiology*, 100:2653 – 2668.
- Goodfellow, L. D. (1938). A psychological interpretation of the results of the zenith radio experiments in telepathy. *Journal of Experimental Psychology*, 23(6):601.
- Green, D. M. and Swets, J. A. (1966). *Signal detection theory and psychophysics*. New York: Wiley.
- Harvey, B. M., Klein, B. P., Petridou, N., and Dumoulin, S. O. (2013). Topographic representation of numerosity in the human parietal cortex. *Science*, 341(6150):1123–1126.
- Helmholtz, H. v. (1867). *Handbuch der Physiologischen Optik*. Allg. Enzyklopädie der Physik 9. Bd. Voss, Leipzig, Germany.
- Hertwig, R. and Erev, I. (2009). The description–experience gap in risky choice. *Trends in cognitive sciences*, 13(12):517–523.
- Hurewitz, F., Gelman, R., and Schnitzer, B. (2006). Sometimes area counts more than number. *Proceedings of the National Academy of Sciences*, 103(51):19599–19604.
- Hürlimann, F., Kiper, D. C., and Carandini, M. (2002). Testing the bayesian model of perceived speed. *Vision research*, 42(19):2253–2257.
- Jacobs, R. A. (1999). Optimal integration of texture and motion cues to depth. *Vision research*, 39(21):3621–3629.
- Jarvstad, A., Hahn, U., Rushton, S. K., and Warren, P. A. (2013). Perceptuo-motor, cognitive, and description-based decision-making seem equally good. *Proceedings of the National Academy of Sciences*, 110(40):16271–16276.
- Jarvstad, A., Rushton, S. K., Warren, P. A., and Hahn, U. (2012). Knowing when to move on: Cognitive and perceptual decisions in time. *Psychological science*, 23(6):589–597.
- Jazayeri, M. and Movshon, J. (2007). A new perceptual illusion reveals mechanisms of sensory decoding. *Nature*, 446:912ff.
- Jazayeri, M. and Shadlen, M. N. (2010). Temporal context calibrates interval timing. *Nature neuroscience*, 13(8):1020.
- Johansson, P., Hall, L., Sikström, S., and Olsson, A. (2005). Failure to detect mismatches between intention and outcome in a simple decision task. *Science*, 310(5745):116–119.

- Jonas, E., Schulz-Hardt, S., Frey, D., and Thelen, N. (2001). Confirmation bias in sequential information search after preliminary decisions: an expansion of dissonance theoretical research on selective exposure to information. *Journal of personality and social psychology*, 80(4):557.
- Kac, M. (1962). A note on learning signal detection. *IRE transactions on information theory*, 8(2):126–128.
- Kahneman, D. (2003). Maps of bounded rationality: Psychology for behavioral economics. *American economic review*, 93(5):1449–1475.
- Kahneman, D. (2011). *Thinking, Fast and Slow*. Macmillan.
- Kahneman, D. and Tversky, A. (1984). Choices, values, and frames. *American Psychologist*, 39(4):341–350.
- Kasparov, G. (2017). *Deep Thinking: Where Machine Intelligence Ends and Human Creativity Begins*. PublicAffairs.
- Kepecs, A., Uchida, N., Zariwala, H. A., and Mainen, Z. F. (2008). Neural correlates, computation and behavioural impact of decision confidence. *Nature*, 455:227 EP –.
- Kersten, D., Mamassian, P., and Yuille, A. (2004). Object perception as Bayesian inference. *Annual Review of Psychology*, 55:271–304.
- Kersten, D. and Schrater, P. (2002). Pattern inference theory: A probabilistic approach to vision. In *Perception and the physical world: Psychological and philosophical issues in perception*. Wiley.
- Kiani, R., Hanks, T. D., and Shadlen, M. N. (2008). Bounded integration in parietal cortex underlies decisions even when viewing duration is dictated by the environment. *Journal of Neuroscience*, 28(12):3017–3029.
- Kim, S. and Burge, J. (2018). The lawful imprecision of human surface tilt estimation in natural scenes. *eLife*, 7:e31448.
- Kim, T. D., Kabir, M., and Gold, J. I. (2017). Coupled decision processes update and maintain saccadic priors in a dynamic environment. *Journal of Neuroscience*, 37(13):3632–3645.
- Kingdom, F. and Prins, N. (2010). *Psychophysics: a practical introduction*. Academic Press London.
- Knill, D. and Richards, W., editors (1996). *Perception as Bayesian Inference*. Cambridge University Press.
- Knill, D. C. (2007). Learning Bayesian priors for depth perception. *Journal of vision*, 7(8):13.
- Knox, R. E. and Inkster, J. A. (1968). Postdecision dissonance at post time. *Journal of personality and social psychology*, 8(4p1):319.

- Kononenko, I. (2001). Machine learning for medical diagnosis: history, state of the art and perspective. *Artificial Intelligence in medicine*, 23(1):89–109.
- Körding, K. and Wolpert, D. (2004). Bayesian integration in sensorimotor learning. *Nature*, 427(15):244–247.
- Körding, K. P., Beierholm, U., Ma, W. J., Quartz, S., Tenenbaum, J. B., and Shams, L. (2007). Causal inference in multisensory perception. *PLoS one*, 2(9):e943.
- Körding, K. P. and Wolpert, D. M. (2006). Bayesian decision theory in sensorimotor control. *Trends in cognitive sciences*, 10(7):319–326.
- Lagnado, D. and Shanks, D. (2003). The influence of hierarchy on probability judgment. *Cognition*, 89(2):157–178.
- Lake, B. M., Ullman, T. D., Tenenbaum, J. B., and Gershman, S. J. (2017). Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40.
- Laming, D. (1979). Autocorrelation of choice-reaction times. *Acta Psychologica*, 43:381–412.
- Landy, M. S., Goutcher, R., Trommershäuser, J., and Mamassian, P. (2007). Visual estimation under risk. *Journal of vision*, 7(6):4–4.
- Lecky, P. (1945). *Self-consistency, a theory of personality*. Island press, New York.
- Lee, S. W. and Schwarz, N. (2010). Washing away postdecisional dissonance. *Science*, 328(5979):709–709.
- Leuthesser, L., Kohli, C. S., and Harich, K. R. (1995). Brand equity: the halo effect measure. *European Journal of Marketing*, 29(4):57–66.
- Link, S. (1975). The relative judgment theory of two choice response time. *Journal of Mathematical Psychology*, 12:114–135.
- Loftus, E. F. (1975). Leading questions and the eyewitness report. *Cognitive psychology*, 7(4):560–572.
- Loftus, E. F. (2005). Planting misinformation in the human mind: A 30-year investigation of the malleability of memory. *Learning & Memory*, 12(4):361–366.
- Loftus, E. F. and Loftus, G. R. (1980). On the permanence of stored information in the human brain. *American Psychologist*, 35(5):409.
- Lu, Z.-L. and Doshier, B. (2013). *Visual psychophysics: From laboratory to theory*. MIT Press.
- Luu, L., Qiu, C., and Stocker, A. (2017). High- to low-level decoding does not generally improve perceptual performance. *bioRxiv*.
- Luu, L. and Stocker, A. (2018). Self-consistent-model. *GitHub*, (2582b6b.).

- Luu, L. and Stocker, A. A. (2016). Making a categorical decision does not modify the stimulus representation in working memory. In *Vision Science Society VSS conference*. Poster presentation.
- Ma, W., Beck, J., Latham, P., and Pouget, A. (2006). Bayesian inference with probabilistic population codes. *Nature Neuroscience*, 9:1432ff.
- Macmillan, N. A. and Creelman, C. D. (2004). *Detection theory: A user's guide*. Psychology press.
- Maloney, L. (2002). Statistical decision theory and biological vision. In *Perception and the physical world: Psychological and philosophical issues in perception*. Wiley.
- Marks, L. E. (1974). On scales of sensation: Prolegomena to any future psychophysics that will be able to come forth as science. *Perception & Psychophysics*, 16(2):358–376.
- Maxwell, J. C. (1857). Xviii.experiments on colour, as perceived by the eye, with remarks on colour-blindness. *Earth and Environmental Science Transactions of the Royal Society of Edinburgh*, 21(2):275–298.
- McGill, W. J. (1957). Serial effects in auditory threshold judgments. *Journal of Experimental Psychology*, 53(5):297.
- Meindertsma, T., Kloosterman, N. A., Nolte, G., Engel, A. K., and Donner, T. H. (2017). Multiple transient signals in human visual cortex associated with an elementary decision. *Journal of Neuroscience*, 37(23):5744–5757.
- Meyniel, F., Schlunegger, D., and Dehaene, S. (2015). The sense of confidence during probabilistic learning: A normative account. *PLoS computational biology*, 11(6):e1004305.
- Mynatt, C. R., Doherty, M. E., and Tweney, R. D. (1977). Confirmation bias in a simulated research environment: An experimental study of scientific inference. *The quarterly journal of experimental psychology*, 29(1):85–95.
- Navajas, J., Bahrami, B., and Latham, P. E. (2016). Post-decisional accounts of biases in confidence. *Current Opinion in Behavioral Sciences*, 11:55–60.
- Nickerson, R. (1998). Confirmation bias: A ubiquitous phenomenon in many guises. *Review of General Psychology*, 2(2):175–220.
- Nieder, A. and Merten, K. (2007). A labeled-line code for small and large numerosities in the monkey prefrontal cortex. *Journal of Neuroscience*, 27(22):5986–5993.
- Nieder, A. and Miller, E. K. (2003). Coding of cognitive magnitude: Compressed scaling of numerical information in the primate prefrontal cortex. *Neuron*, 37(1):149–157.
- Nienborg, H. and Cumming, B. (2009). Decision-related activity in sensory neurons reflects more than a neuron's causal effect. *Nature*, 459:89–92.
- Nisbett, R. E. and Wilson, T. D. (1977). The halo effect: Evidence for unconscious alteration of judgments. *Journal of personality and social psychology*, 35(4):250.

- Oaksford, M. and Hall, S. (2016). On the source of human irrationality. *Trends in cognitive sciences*, 20(5):336–344.
- Obermeyer, Z. and Emanuel, E. J. (2016). Predicting the futurebig data, machine learning, and clinical medicine. *The New England journal of medicine*, 375(13):1216.
- O’Callaghan, C., Kveraga, K., Shine, J. M., Adams, R. B., and Bar, M. (2016). Convergent evidence for top-down effects from the predictive brain 1. *Behavioral and Brain Sciences*, 39.
- Peters, M., Thesen, T., Ko, Y., Maniscalco, B., Carslon, C., Davidson, M., Doyle, W., Kuzniecky, R., Devinsky, O., Halgren, E., and Lau, H. (2017). Perceptual confidence neglects decision-incongruent evidence in the brain. *Nature Human Behaviour*, 1:1–8.
- Piazza, M., Izard, V., Pinel, P., Le Bihan, D., and Dehaene, S. (2004). Tuning curves for approximate numerosity in the human intraparietal sulcus. *Neuron*, 44(3):547–555.
- Pitkow, X. and Angelaki, D. (2017). Inference in the brain: Statistics flowing in redundant population codes. *Neuron*, 94:943–953.
- Pleskac, T. J. and Busemeyer, J. R. (2010). Two-stage dynamic signal detection: a theory of choice, decision time, and confidence. *Psychological review*, 117(3):864.
- Preston, M. G. (1936). Contrast effects and the psychophysical judgments. *The American Journal of Psychology*, 48(3):389–402.
- Rahnev, D., Koizumi, A., McCurdy, L., D’Esposito, M., and Lau, H. (2015). Confidence leak in perceptual decision making. *Psychological Science*, 26(11):1664–1680.
- Rakow, T. and Newell, B. R. (2010). Degrees of uncertainty: An overview and framework for future research on experience-based choice. *Journal of Behavioral Decision Making*, 23(1):1–14.
- Riefer, P. S., Prior, R., Blair, N., Pavey, G., and Love, B. C. (2017). Coherency-maximizing exploration in the supermarket. *Nature human behaviour*, 1(1):0017.
- Sato, Y. and Kording, K. P. (2014). How much to trust the senses: Likelihood learning. *Journal of Vision*, 14(13):13–13.
- Senders, V. L. (1953). Further analysis of response sequences in the setting of a psychophysical experiment. *The American journal of psychology*, 66(2):215–228.
- Senders, V. L. and Sowards, A. (1952). Analysis of response sequences in the setting of a psychophysical experiment. *The American journal of psychology*, 65(3):358–374.
- Seriès, P. and Seitz, A. (2013). Learning what to expect (in visual perception). *Frontiers in human neuroscience*, 7:668.
- Sharot, T., Fleming, S. M., Yu, X., Koster, R., and Dolan, R. J. (2012). Is choice-induced preference change long lasting? *Psychological science*, 23(10):1123–1129.

- Sharot, T., Velasquez, C., and Dolan, R. (2010). Do decisions shape preference? Evidence from blind choice. *Psychological Science*, 21(9):1231–1235.
- Shi, Z., Church, R. M., and Meck, W. H. (2013). Bayesian optimization of time perception. *Trends in Cognitive Sciences*, 17(11):556–564.
- Siegel, M., Buschman, T., and Miller, E. (2015). Cortical information flow during flexible sensorimotor decisions. *Science*, 348(6241):1352–1355.
- Sigall, H. and Ostrove, N. (1975). Beautiful but dangerous: effects of offender attractiveness and nature of the crime on juridic judgment. *Journal of Personality and Social Psychology*, 31(3):410.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. (2017). Mastering the game of go without human knowledge. *Nature*, 550(7676):354.
- Simon, D., Snow, C. J., and Read, S. J. (2004). The redux of cognitive consistency theories: evidence judgments by constraint satisfaction. *Journal of personality and social psychology*, 86(6):814.
- Simoncelli, E. (2009). *The New Cognitive Neuroscience*, chapter Optimal Estimation in Sensory Systems. MIT Press, 4th edition.
- Singer, P. W. (2010). War of the machines. *Scientific American*, 303(1):56–63.
- Staw, B. M. (1976). Knee-deep in the big muddy: A study of escalating commitment to a chosen course of action. *Organizational behavior and human performance*, 16(1):27–44.
- Stevens, S. (1946). On the theory of scales of measurement. *Science*, 103(2684):677–680.
- Stevens, S. S. (1956). The direct estimation of sensory magnitudes: Loudness. *The American journal of psychology*, 69(1):1–25.
- Stevens, S. S. (1957). On the psychophysical law. *Psychological review*, 64(3):153.
- Stocker, A. and Simoncelli, E. (2006). Noise characteristics and prior expectations in human visual speed perception. *Nature Neuroscience*, pages 578–585.
- Stocker, A. and Simoncelli, E. (2007). A Bayesian model of conditioned perception. In Platt, J., Koller, D., Singer, Y., and Roweis, S., editors, *Advances in Neural Information Processing Systems NIPS 20*, pages 1409–1416, Cambridge, MA. MIT Press.
- Szpiro, S. F. A., Spering, M., and Carrasco, M. (2014). Perceptual learning modifies untrained pursuit eye movements. *Journal of Vision*, 14(8):8.
- Tank, A. and Stocker, A. A. (2014). Biased perception leads to biased action: Validating a bayesian model of interception. *arXiv preprint arXiv:1403.1920*.

- Thorndike, E. (1920). A constant error in psychological ratings. *Journal of Applied Psychology*, 4(1):25.
- Treisman, M. and Williams, T. C. (1984). A theory of criterion setting with an application to sequential dependencies. *Psychological Review*, 91(1):68.
- Trommershäuser, J., Maloney, L. T., and Landy, M. S. (2003). Statistical decision theory and the selection of rapid, goal-directed movements. *JOSA A*, 20(7):1419–1433.
- Trommershäuser, J., Maloney, L. T., and Landy, M. S. (2008). Decision making, movement planning and statistical decision theory. *Trends in cognitive sciences*, 12(8):291–297.
- Tsetsos, K., Gao, J., McClelland, J. L., and Usher, M. (2012). Using time-varying evidence to test models of decision dynamics: bounded diffusion vs. the leaky competing accumulator model. *Frontiers in neuroscience*, 6:79.
- Tsetsos, K., Moran, R., Moreland, J., Chater, N., Usher, M., and Summerfield, C. (2016). Economic irrationality is optimal during noisy decision making. *Proceedings of the National Academy of Sciences*, 113(11):3102–3107.
- Tune, G. (1964). Response preferences: A review of some relevant literature. *Psychological bulletin*, 61(4):286.
- Tversky, A. and Kahneman, D. (1974). Judgment under uncertainty: Heuristics and biases. *science*, 185(4157):1124–1131.
- Tversky, A. and Kahneman, D. (1981). The framing of decisions and the psychology of choice. *science*, 211(4481):453–458.
- Urai, A. E., Braun, A., and Donner, T. H. (2017). Pupil-linked arousal is driven by decision uncertainty and alters serial choice bias. *Nature communications*, 8:14637.
- van den Berg, R., Anandalingam, K., Zylberberg, A., Kiani, R., Shadlen, M. N., and Wolpert, D. M. (2016). A common mechanism underlies changes of mind about decisions and confidence. *eLife*, 5:e12192.
- Verplanck, W. S. and Blough, D. S. (1958). Randomized stimuli and the non-independence of successive responses at the visual threshold. *The Journal of general psychology*, 59(2):263–272.
- Verplanck, W. S., Collier, G. H., and Cotton, J. W. (1952). Nonindependence of successive responses in measurements of the visual threshold. *Journal of experimental psychology*, 44(4):273.
- Verplanck, W. S., Cotton, J. W., and Collier, G. H. (1953). Previous training as a determinant of response dependency at the threshold. *Journal of Experimental Psychology*, 46(1):10.
- Vervaeck, K. and Boer, L. (1980). Sequential effects in two-choice reaction time: subjective expectancy and automatic after-effect at short response-stimulus intervals. *Acta Psychologica*, 44:175–190.

- Wandell, B. A. (1995). *Foundations of vision*. Sunderland, MA: Sinauer Associates.
- Weber, E. H. (1867). *De pulsu, resorptione, auditu et tactu*. Leipzig: Koehler.
- Wei, X.-X. and Stocker, A. (2012). Bayesian inference with efficient neural population codes. In *Lecture Notes in Computer Science, Artificial Neural Networks and Machine Learning - ICANN 2012, Lausanne, Switzerland*, volume 7552, pages 523–530, Lausanne, Switzerland. Selected talk presentation.
- Wei, X.-X. and Stocker, A. A. (2015). A Bayesian observer model constrained by efficient coding can explain "anti-Bayesian" percepts. *Nature Neuroscience*, 18(10):1509–1517.
- Wei, X.-X. and Stocker, A. A. (2017). Lawful relation between perceptual bias and discriminability. *Proc. National Academies of Sciences U.S.A.*, 114(38):10244–10249.
- Weiss, Y., Simoncelli, E., and Adelson, E. (2002). Motion illusions as optimal percept. *Nature Neuroscience*, 5(6):598–604.
- Whiteley, L. and Sahani, M. (2008). Implicit knowledge of visual uncertainty guides decisions with asymmetric outcomes. *Journal of vision*, 8(3):2–2.
- Wilken, P. and Ma, W. (2004). A detection theory account of change detection. *Journal of Vision*, 4:1120–1135.
- Wixted, J. T., Mickes, L., Clark, S. E., Gronlund, S. D., and Roediger III, H. L. (2015). Initial eyewitness confidence reliably predicts eyewitness identification accuracy. *American Psychologist*, 70(6):515.
- Wolfe, J. M., Kluender, K. R., Levi, D. M., Bartoshuk, L. M., Herz, R. S., Klatzky, R. L., Lederman, S. J., and Merfeld, D. M. (2006). *Sensation & perception*. Sinauer Sunderland, MA.
- Younger, J. C., Walker, L., and Arrowood, A. J. (1977). Postdecision dissonance at the fair. *Personality and Social Psychology Bulletin*, 3(2):284–287.
- Yu, S., Pleskac, T. J., and Zeigenfuse, M. D. (2015a). Dynamics of postdecisional processing of confidence. *Journal of Experimental Psychology: General*, 144(2):489.
- Yu, S., Pleskac, T. J., and Zeigenfuse, M. D. (2015b). Dynamics of postdecisional processing of confidence. *Journal of Experimental Psychology: General*, 144(2):489.
- Zamboni, E., Ledgeway, T., McGraw, P., and Schluppeck, D. (2016). Do perceptual biases emerge early or late in visual processing? decision-biases in motion perception. *Proc. of Royal Society of London B*, 283(1833).
- Zhang, W. and Luck, S. J. (2008). Discrete fixed-resolution representations in visual working memory. *Nature*, 453(7192):233.
- Zylberberg, A., Barttfeld, P., and Sigman, M. (2012). The construction of confidence in a perceptual decision. *Frontiers in integrative neuroscience*, 6:79.